

Sustainability in AI

Roy Schwartz

Hebrew University of Jerusalem

Guest lecture, Machine Learning CS-433
December 2022



1

THE HEBREW
UNIVERSITY
OF JERUSALEM



A Little about Me

Roy Schwartz



- A senior lecturer at the School of CS at the Hebrew U. of Jerusalem
 - I was a postdoc and a research scientist at The University of Washington and the Allen Institute for AI (AI2)

A Little about Me

Roy Schwartz



- A senior lecturer at the School of CS at the Hebrew U. of Jerusalem
 - I was a postdoc and a research scientist at The University of Washington and the Allen Institute for AI (AI2)
- I study Artificial Intelligence (AI) and focus on Natural Language Processing (NLP)
 - **Understanding** AI models
 - Revealing **biases** in datasets
 - Making AI more **sustainable**

A Little about Me

Roy Schwartz



- A senior lecturer at the School of CS at the Hebrew U. of Jerusalem
 - I was a postdoc and a research scientist at The University of Washington and the Allen Institute for AI (AI2)
- I study Artificial Intelligence (AI) and focus on Natural Language Processing (NLP)
 - **Understanding** AI models
 - Revealing **biases** in datasets
 - Making AI more **sustainable**

A Little about Me

Roy Schwartz



- A senior lecturer at the School of CS at the Hebrew U. of Jerusalem
 - I was a postdoc and a research scientist at The University of Washington and the Allen Institute for AI (AI2)
- I study Artificial Intelligence (AI) and focus on Natural Language Processing (NLP)
 - Understanding AI models
 - Revealing **biases** in datasets
 - Making AI more **sustainable**

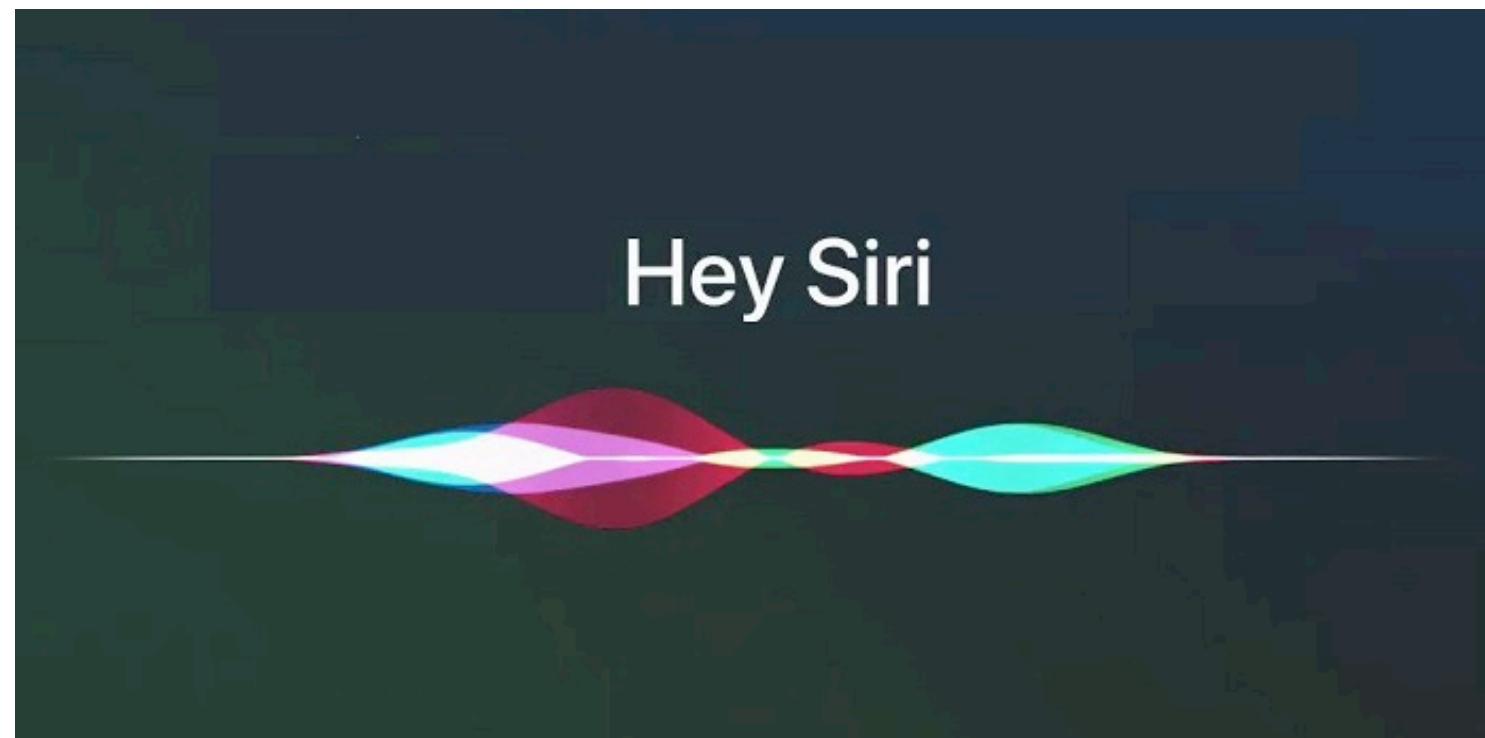
AI Today



Translator

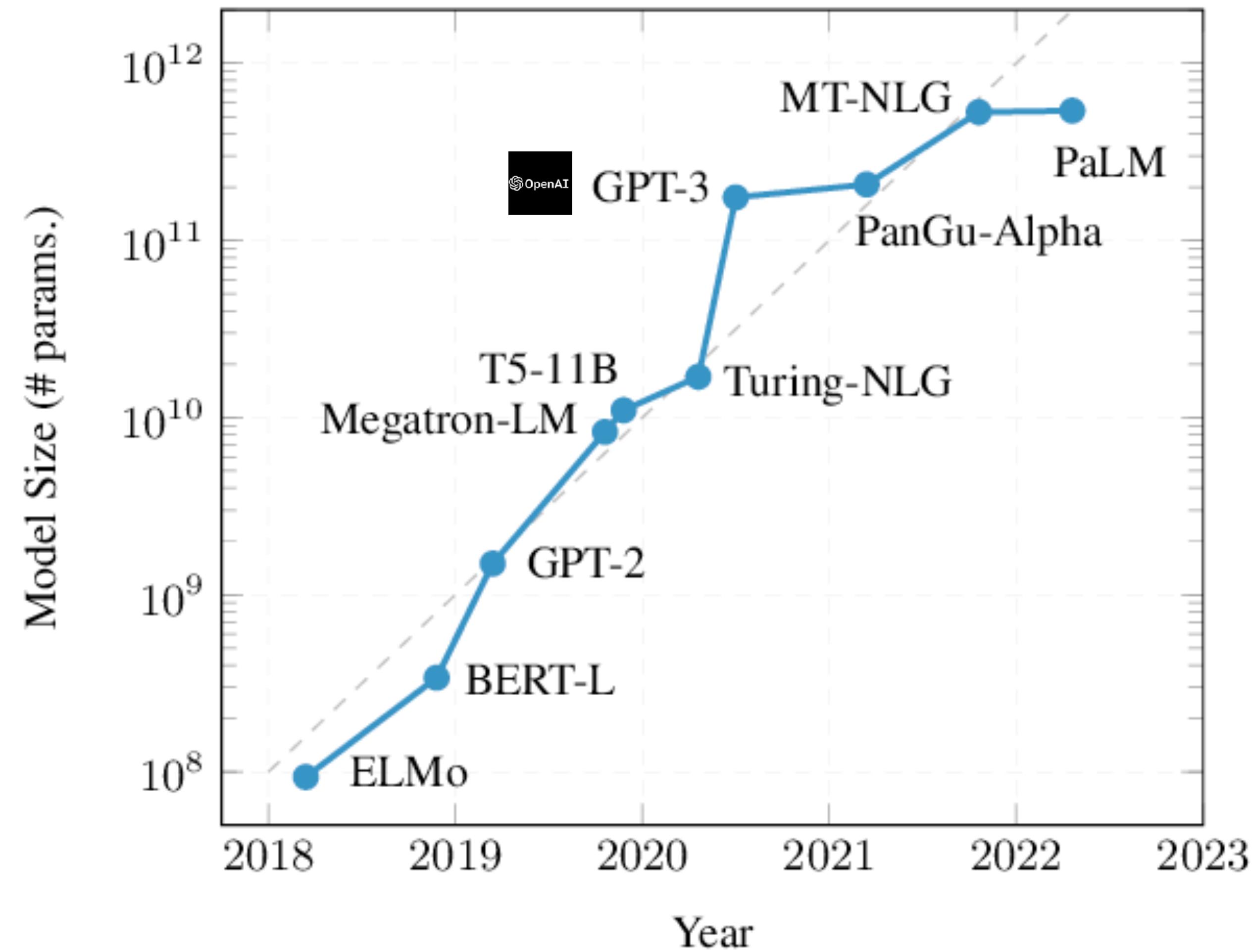


ChatGPT		
Examples	Capabilities	Limitations
"Explain quantum computing in simple terms"	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?"	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?"	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021



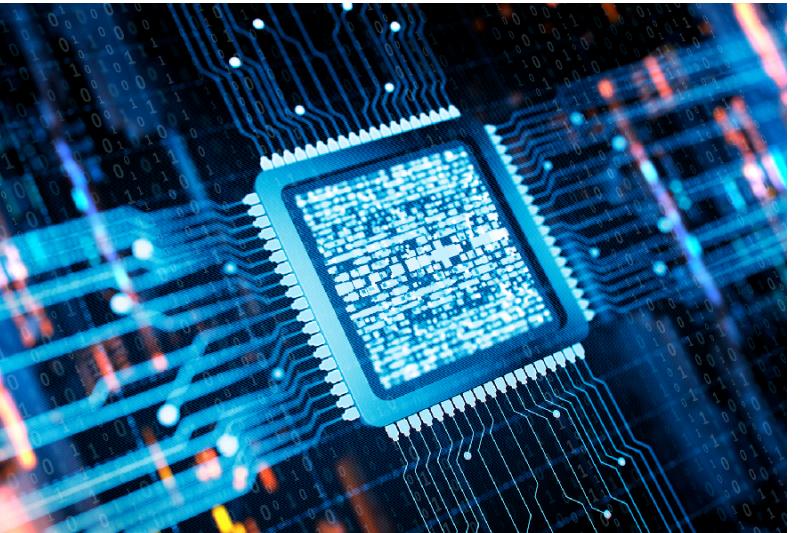
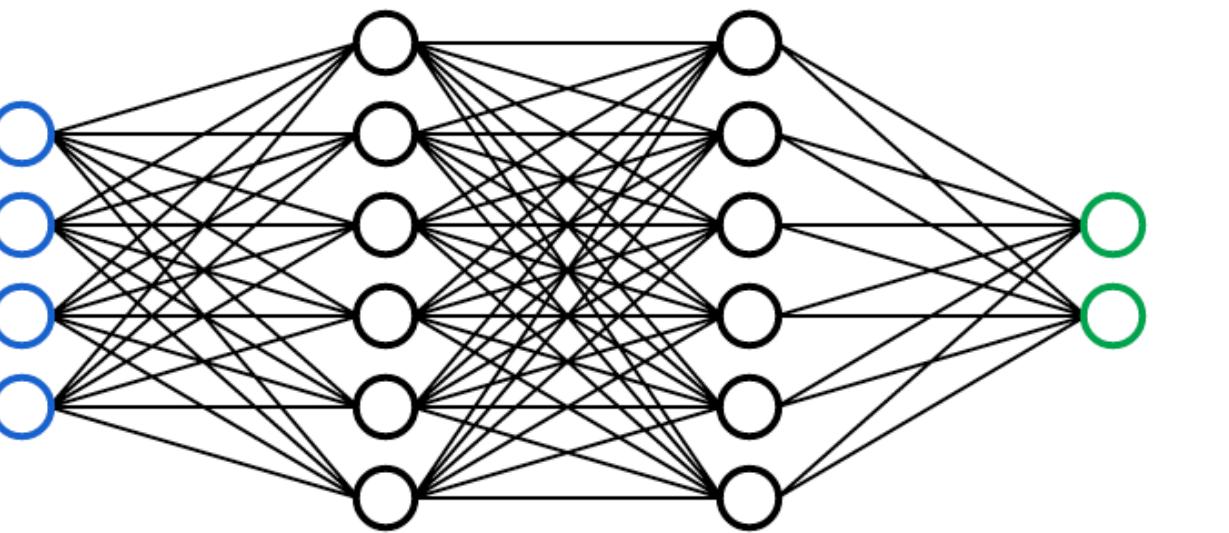
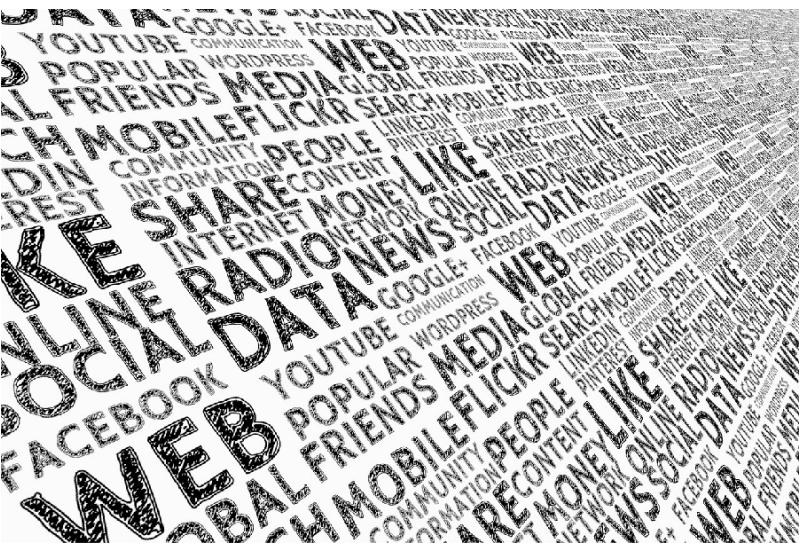
Large Models

5,000X in 4 Years

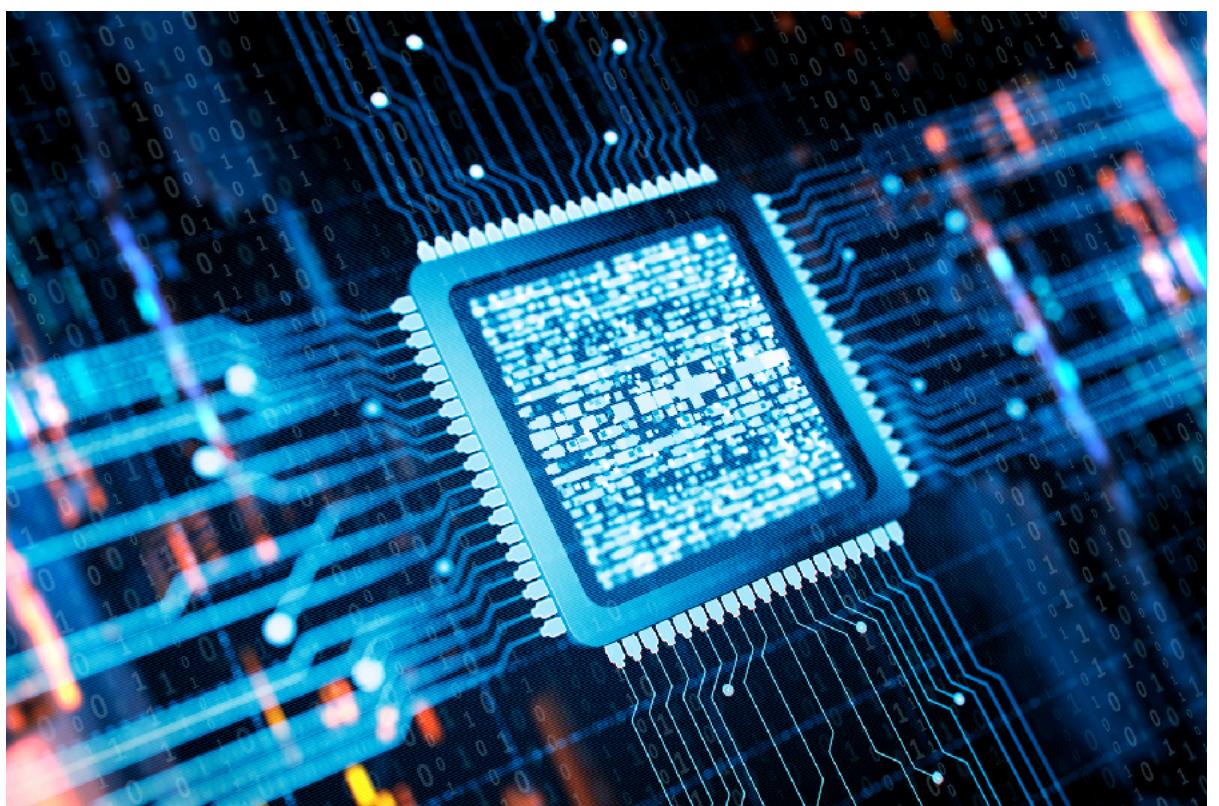
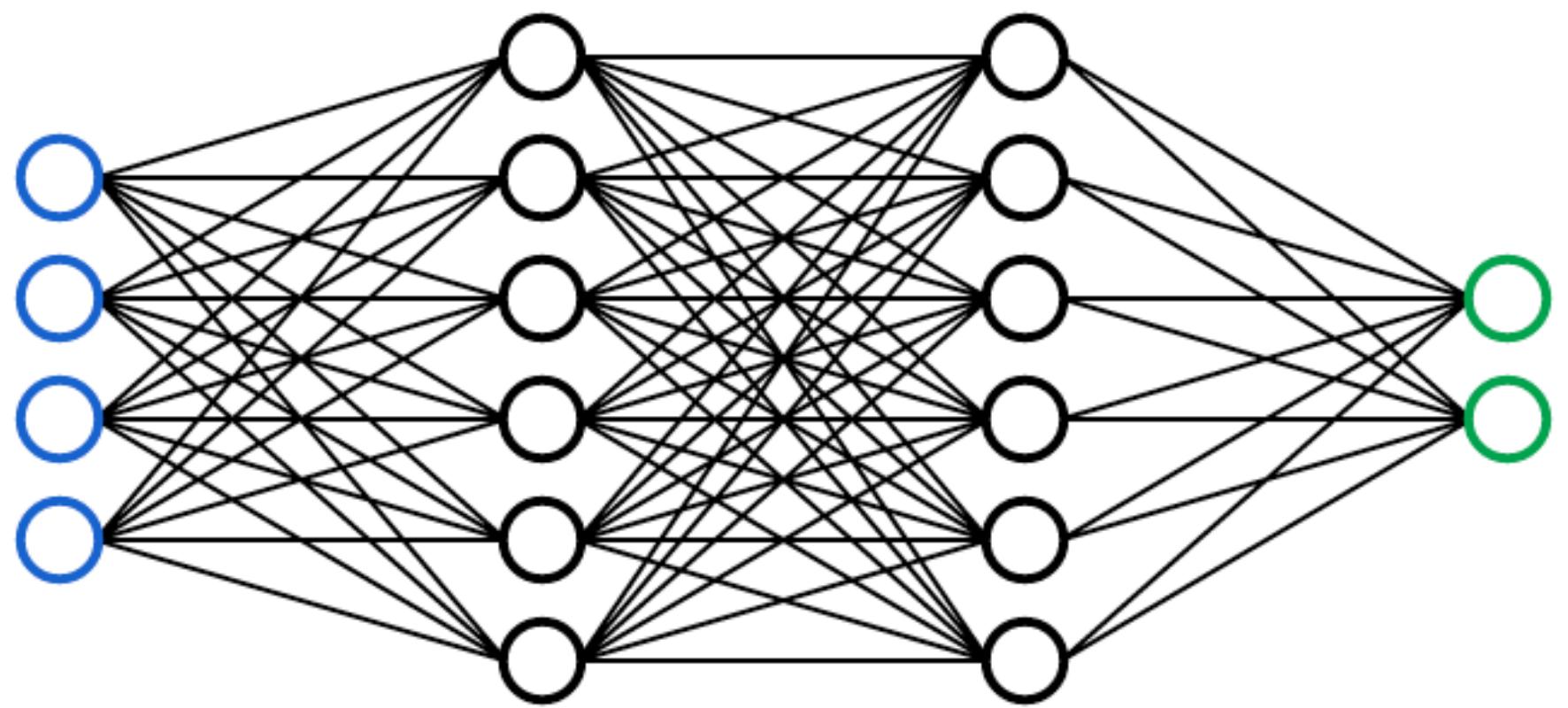


Taken from Lakim et al. (2022)

Scaling Neural Networks

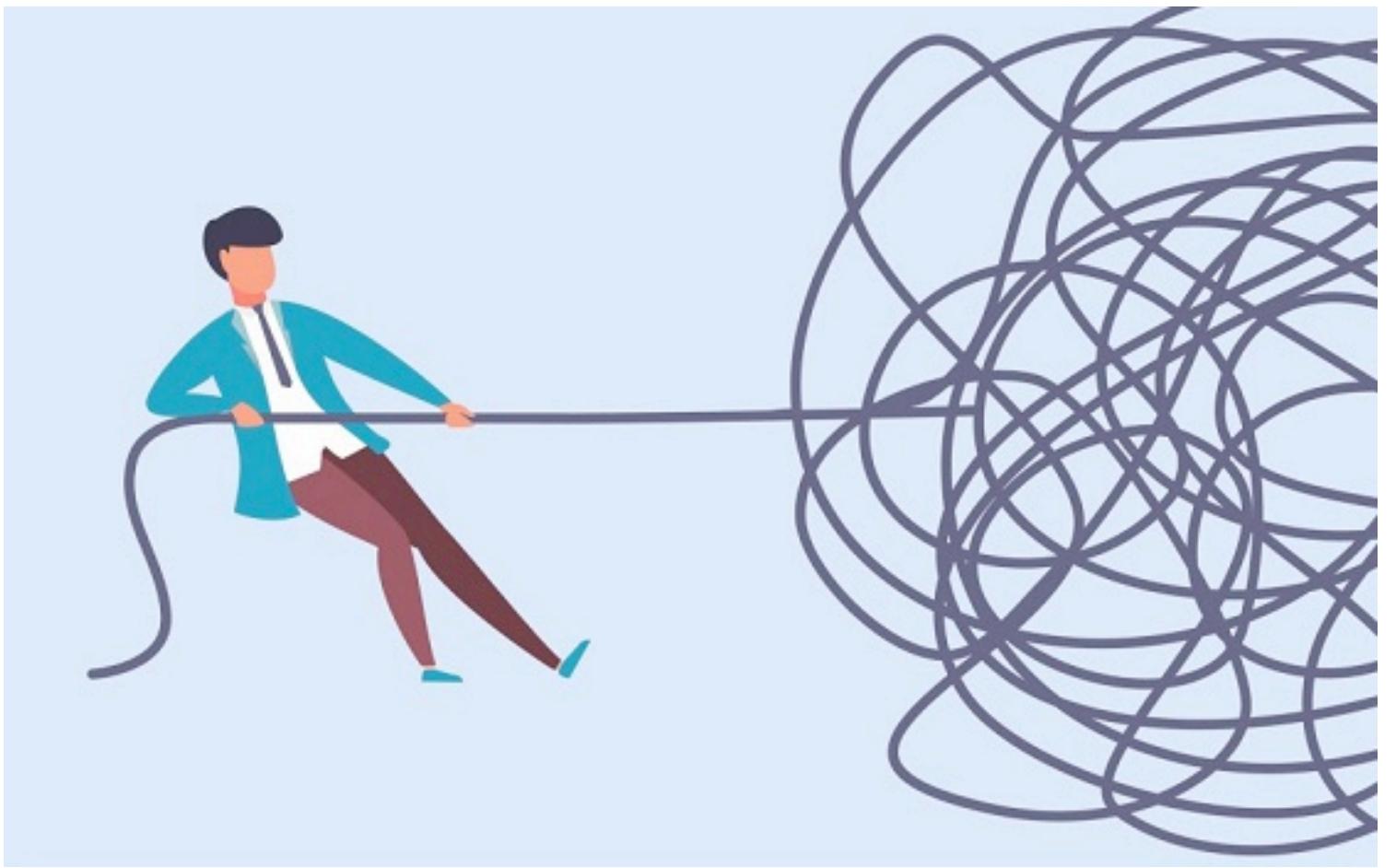


Scaling Neural Networks



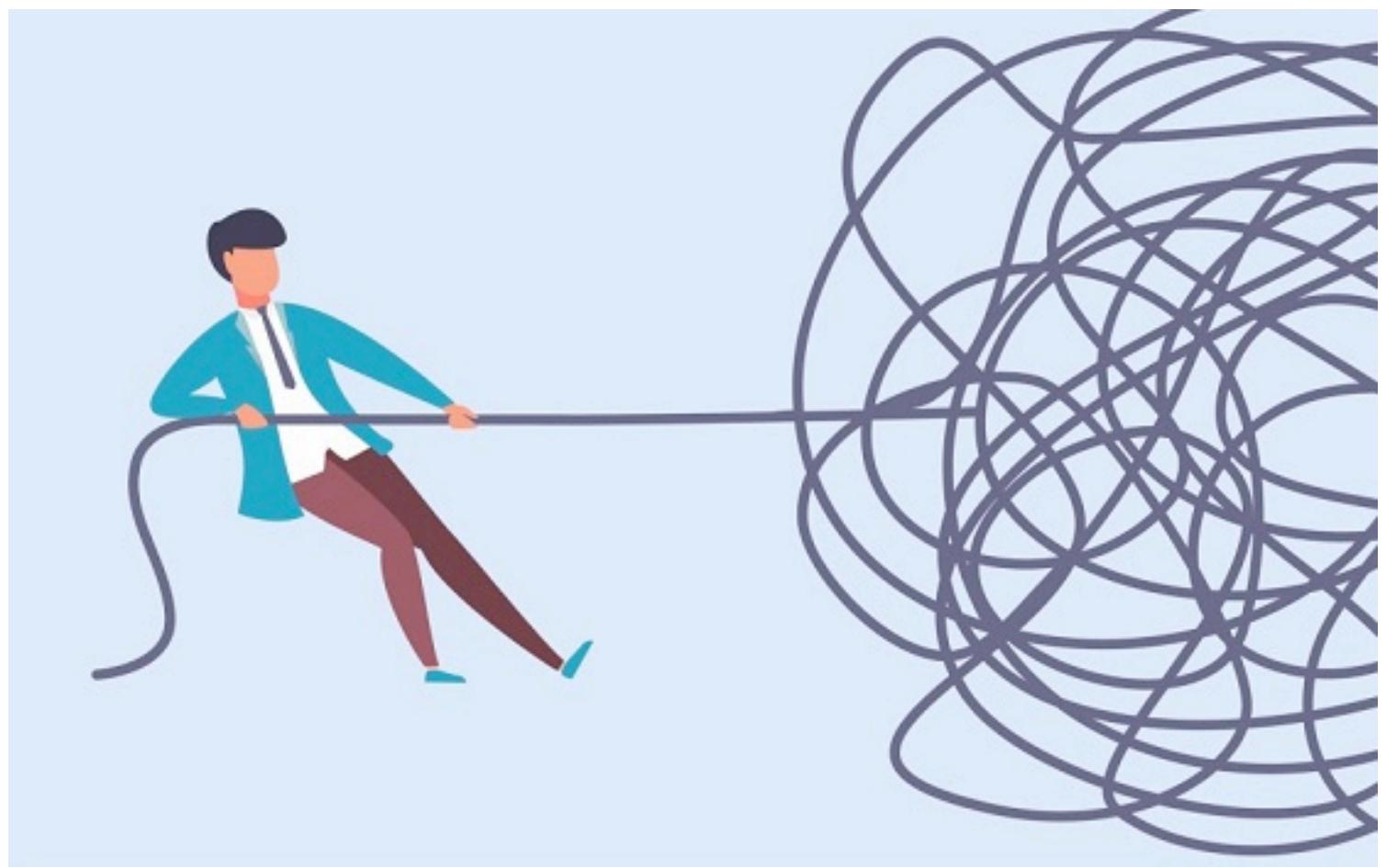
Outline

Outline

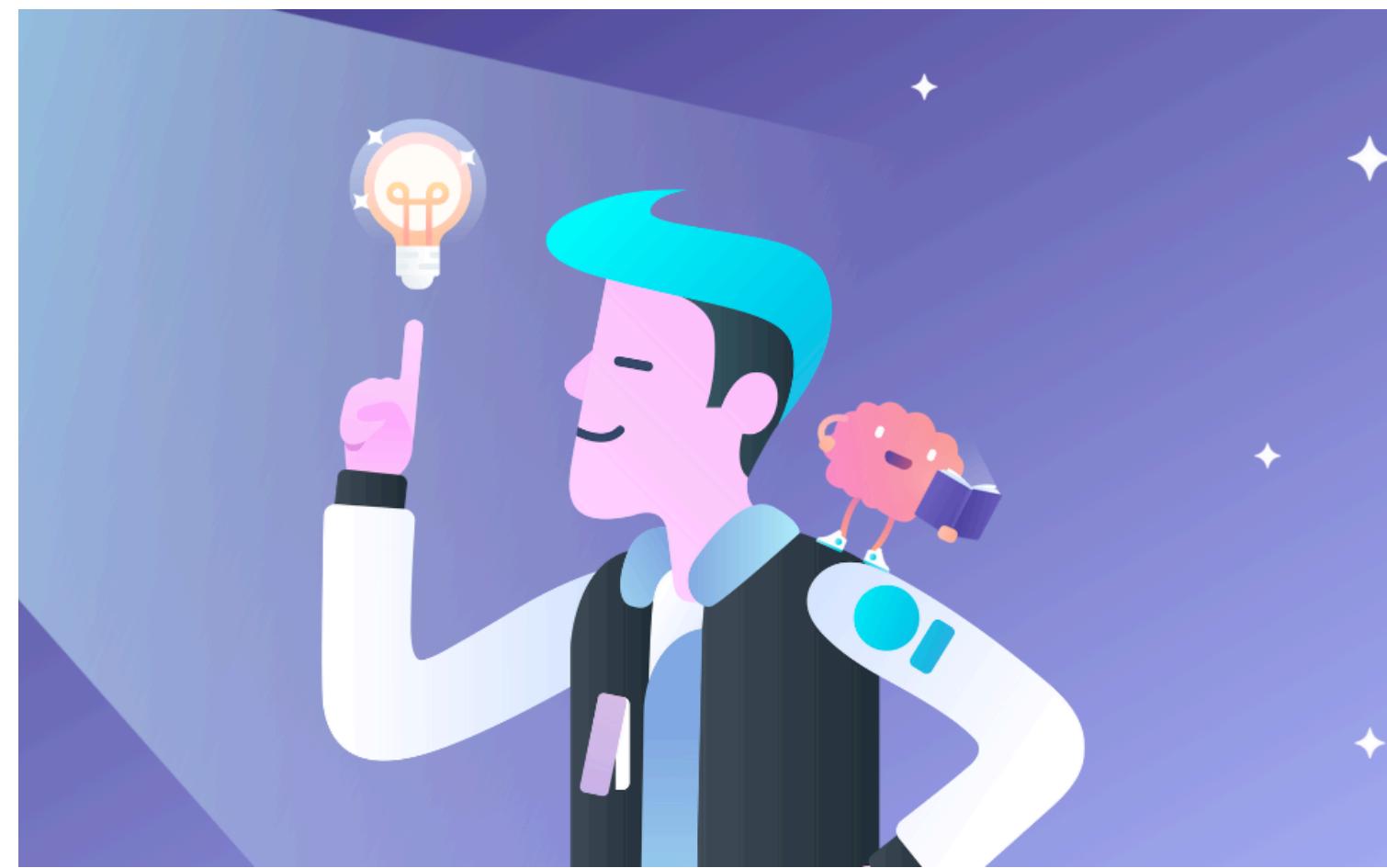


Challenges

Outline

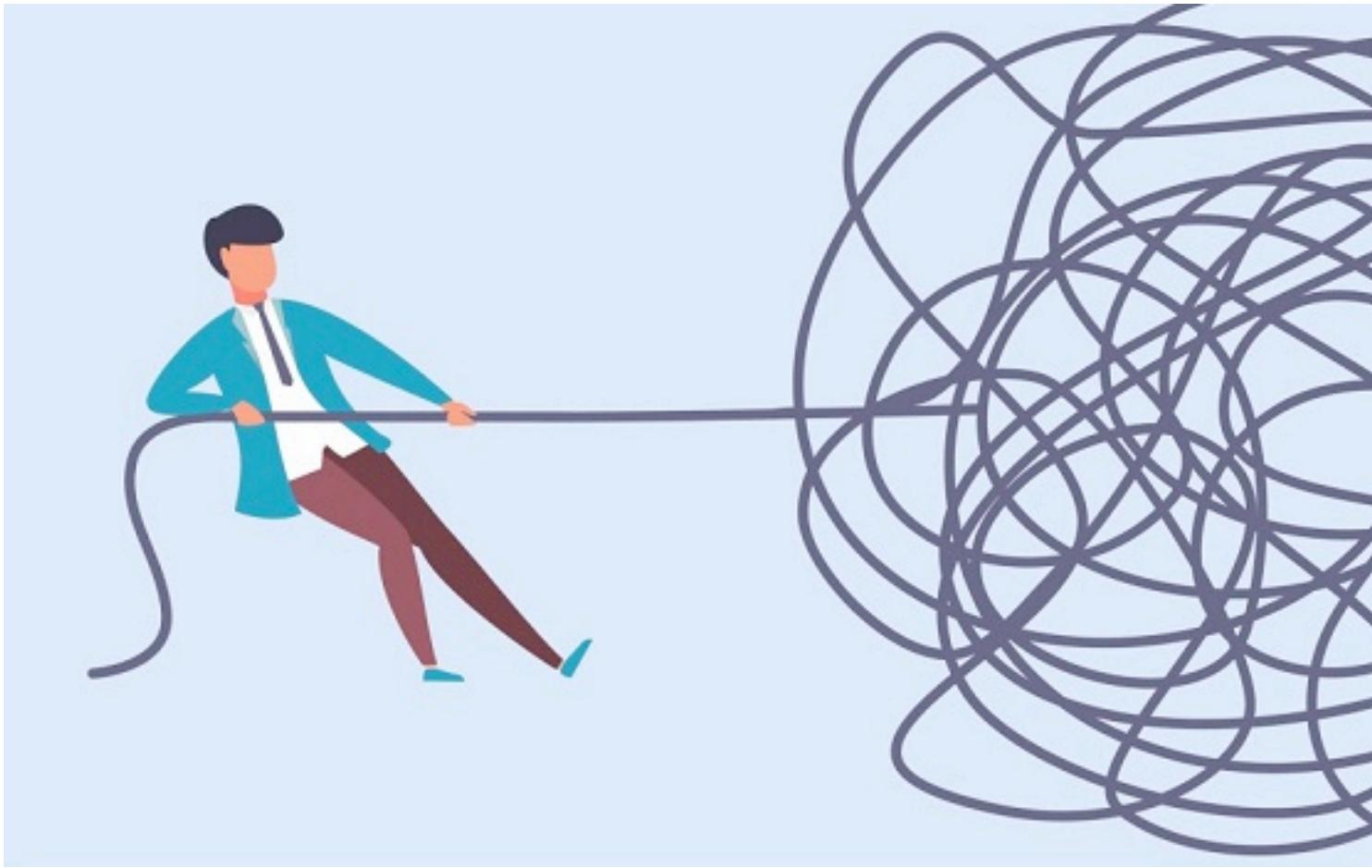


Challenges

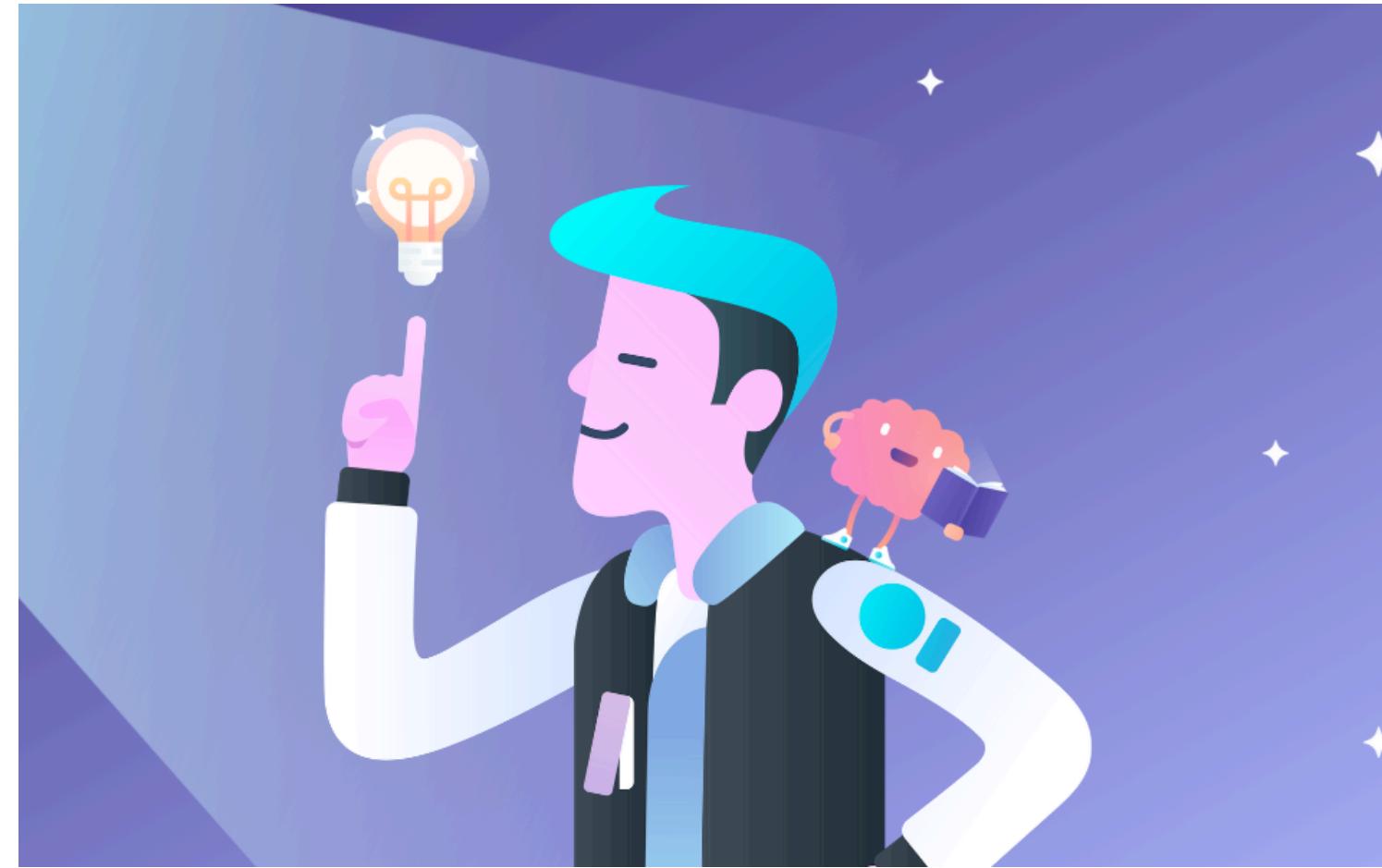


Mitigation

Outline



Challenges

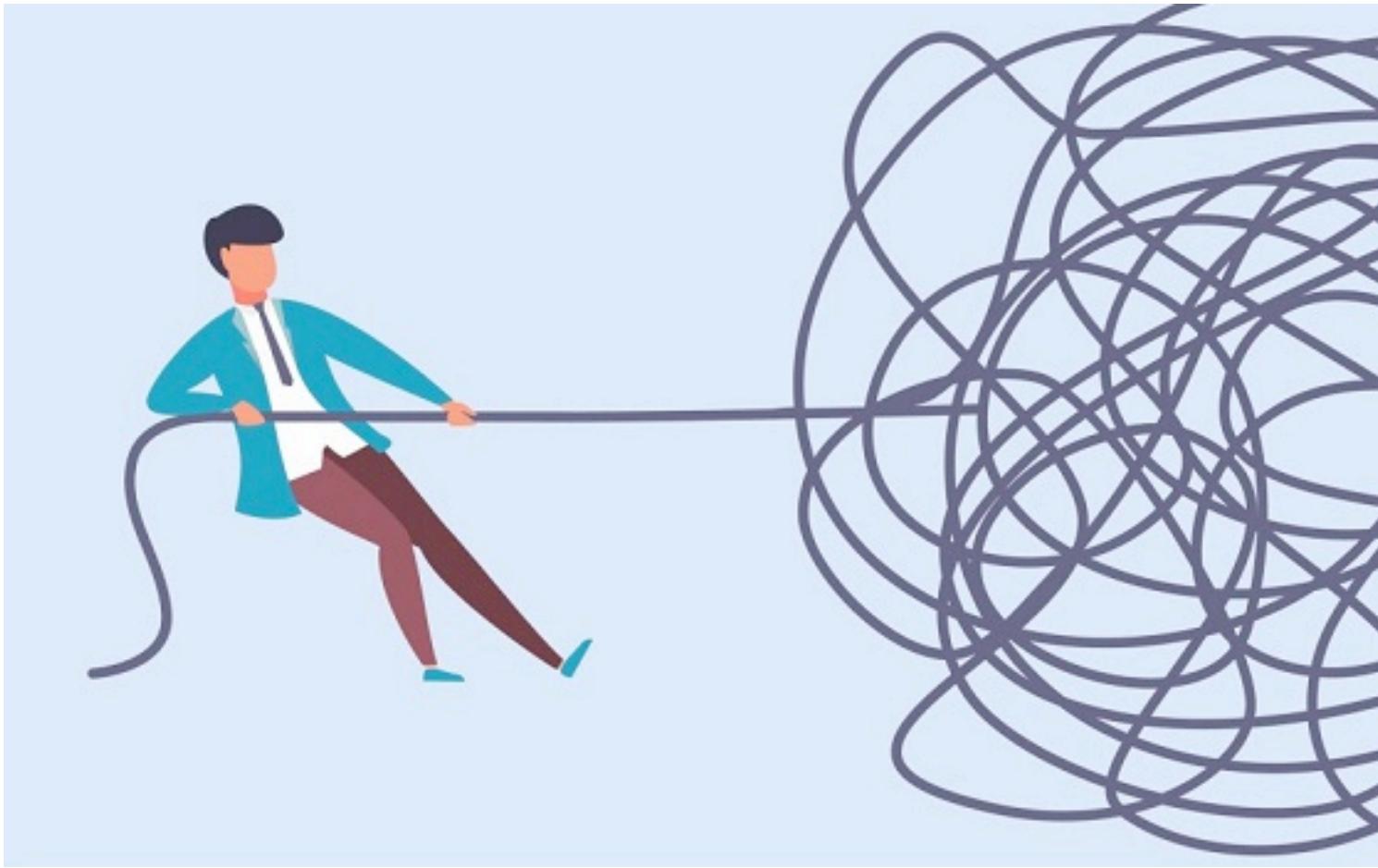


Opportunities

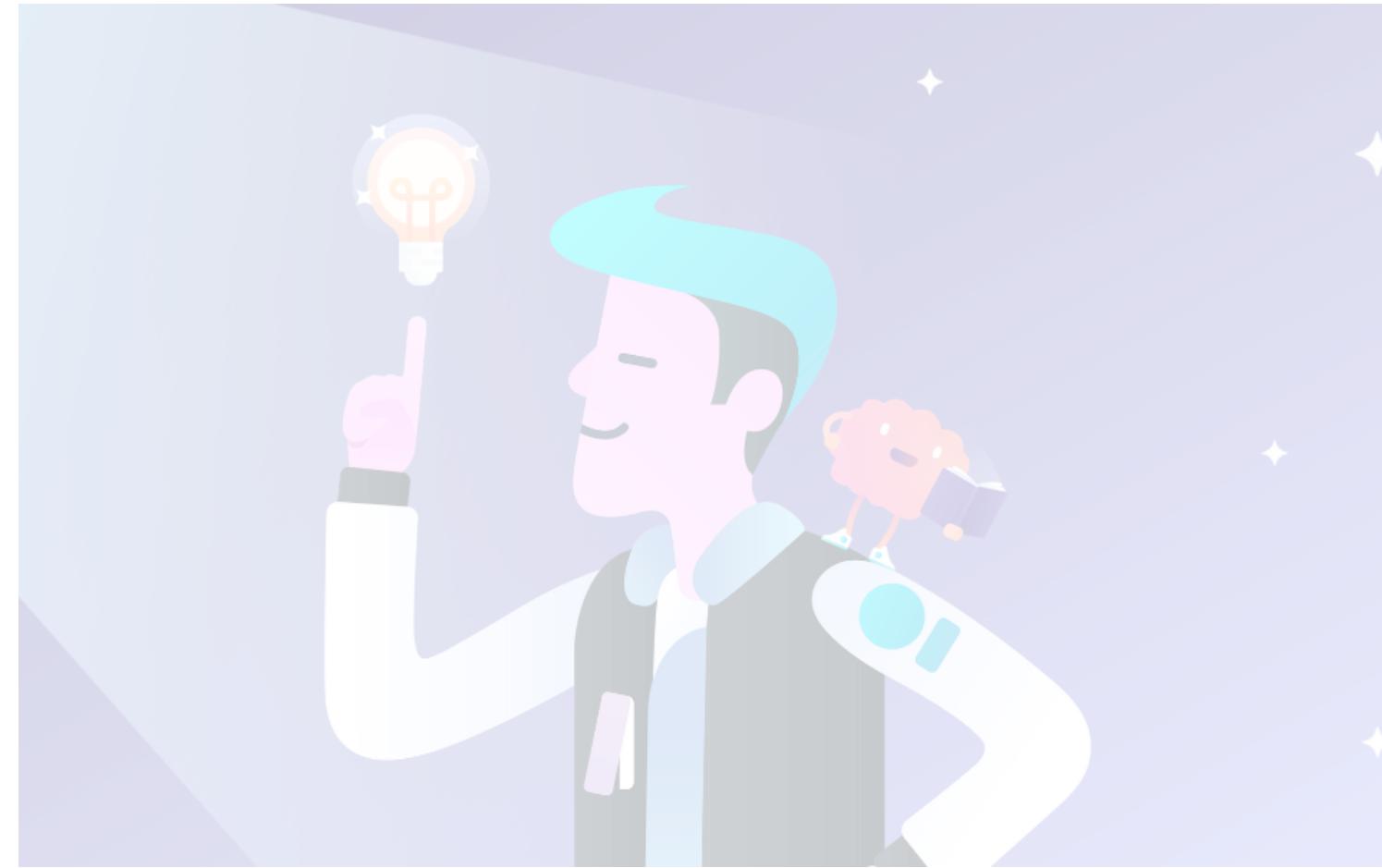


Mitigation

Outline



Challenges



Opportunities



Mitigation



Challenges with Scaling Inclusiveness

Synced
AI TECHNOLOGY & INDUSTRY REVIEW

FEATURE ▾ INDUSTRY ▾ TECHNOLOGY COMMUNITY ▾ ABOUT US ▾ REPORT CONTRIBUTE TO SYNCED REVIEW

The Staggering Cost of Training SOTA AI Models

While it is exhilarating to see AI researchers pushing the performance of cutting-edge models to new heights, the costs of such processes are also rising at a dizzying rate.

<https://syncedreview.com/2019/06/27/the-staggering-cost-of-training-sota-ai-models/>



Training Costs

- BERT (Devlin et al, 2019) was trained on **16** Cloud TPUs for **4** days
- RoBERTa (Liu et al., 2019) was trained on **1024** V100 GPUs for approximately **1** day
- PaLM (Chowdhery et al., 2022) was trained on **6144** TPU v4 chips for **50** days and **3072** TPU v4 chips for **15** days

Number of Authors





Number of Authors

Language Models are Few-Shot Learners

Tom B. Brown* **Benjamin Mann*** **Nick Ryder*** **Melanie Subbiah***

Jared Kaplan† **Prafulla Dhariwal** **Arvind Neelakantan** **Pranav Shyam** **Girish Sastry**

Amanda Askell **Sandhini Agarwal** **Ariel Herbert-Voss** **Gretchen Krueger** **Tom Henighan**

Rewon Child **Aditya Ramesh** **Daniel M. Ziegler** **Jeffrey Wu** **Clemens Winter**

Christopher Hesse **Mark Chen** **Eric Sigler** **Mateusz Litwin** **Scott Gray**

Benjamin Chess **Jack Clark** **Christopher Berner**

Sam McCandlish **Alec Radford** **Ilya Sutskever** **Dario Amodei**

OpenAI



Number of Authors

Language Models are Few-Shot Learners

Tom B. Brown* Benjamin Mann* Nick Ryder* Melanie Subbiah*
Jared Kaplan† Prafulla Dhariwal Arvind Neelakantan Pranav Shyam Girish Sastry
Amanda Askell Sandhini Agarwal Ariel Herbert-Voss Gretchen Krueger Tom Henighan
Rewon Child Aditya Ramesh Daniel M. Ziegler Jeffrey Wu Clemens Winter
Christopher Hesse Mark Chen Eric Sigler Mateusz Litwin Scott Gray
Benjamin Chess Jack Clark Christopher Berner
Sam McCandlish Alec Radford Ilya Sutskever Dario Amodei

OpenAI

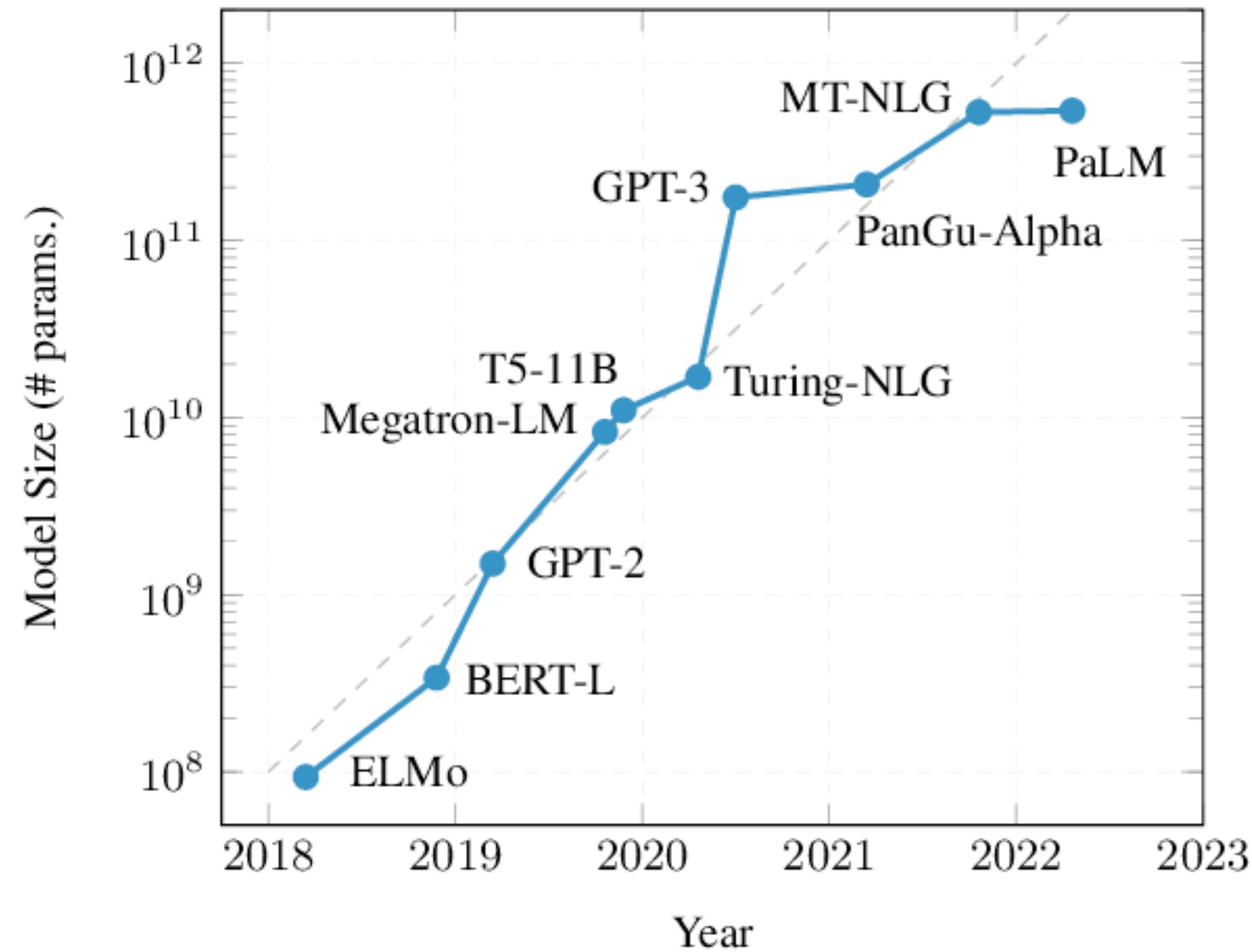
PaLM: Scaling Language Modeling with Pathways

Aakanksha Chowdhery* Sharan Narang* Jacob Devlin*
Maarten Bosma Gaurav Mishra Adam Roberts Paul Barham
Hyung Won Chung Charles Sutton Sebastian Gehrmann Parker Schuh Kensen Shi
Sasha Tsvyashchenko Joshua Maynez Abhishek Rao† Parker Barnes Yi Tay
Noam Shazeer† Vinodkumar Prabhakaran Emily Reif Nan Du Ben Hutchinson
Reiner Pope James Bradbury Jacob Austin Michael Isard Guy Gur-Ari
Pengcheng Yin Toju Duke Anselm Levskaya Sanjay Ghemawat Sunipa Dev
Henryk Michalewski Xavier Garcia Vedant Misra Kevin Robinson Liam Fedus
Denny Zhou Daphne Ippolito David Luan† Hyeontaek Lim Barret Zoph
Alexander Spiridonov Ryan Sepassi David Dohan Shivani Agrawal Mark Omernick
Andrew M. Dai Thanumalayan Sankaranarayana Pillai Marie Pellat Aitor Lewkowycz
Erica Moreira Rewon Child Oleksandr Polozov† Katherine Lee Zongwei Zhou
Xuezhi Wang Brennan Saeta Mark Diaz Orhan Firat Michele Catasta† Jason Wei
Kathy Meier-Hellstern Douglas Eck Jeff Dean Slav Petrov Noah Fiedel

Google Research

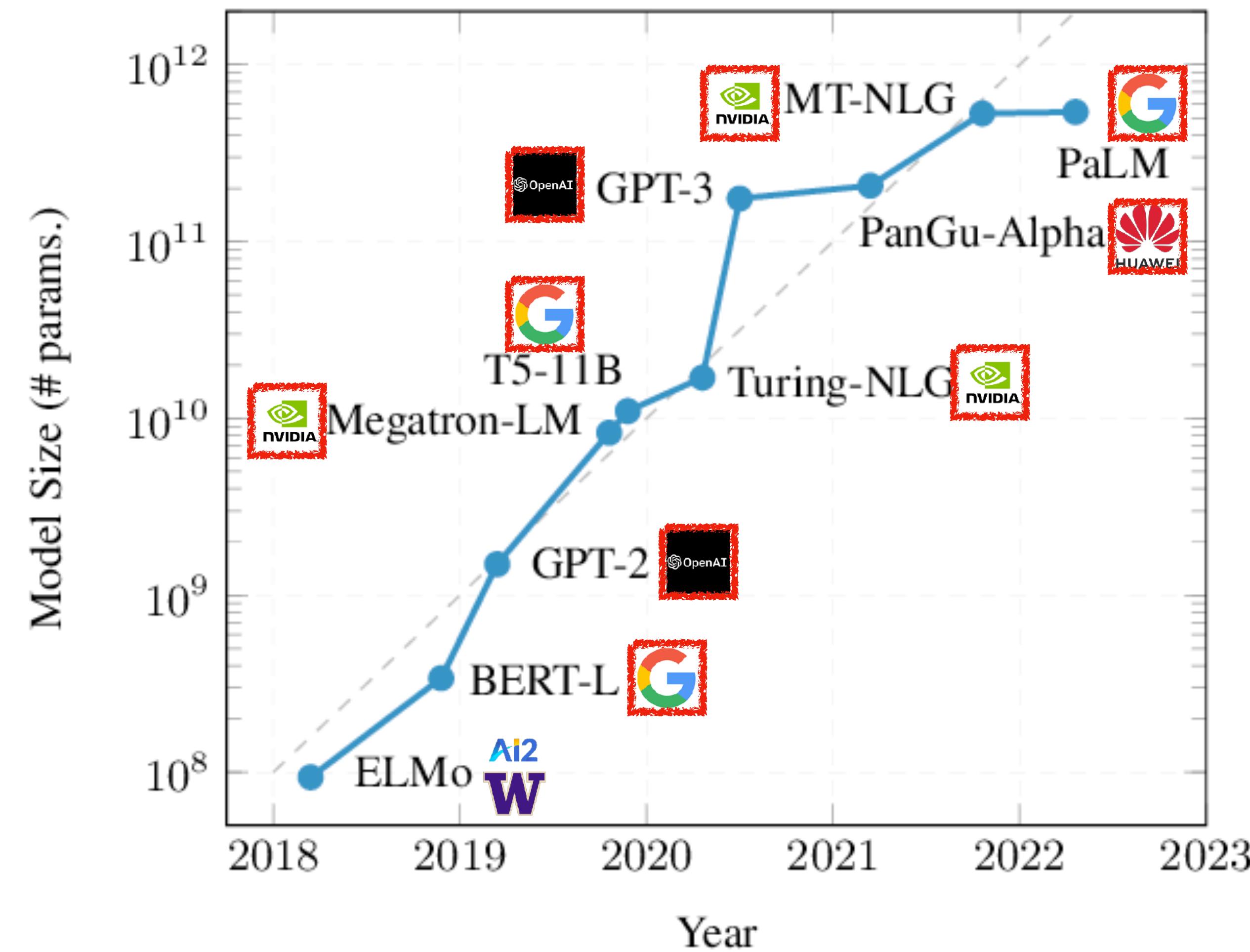


It's a Rich Man's World





It's a Rich Man's World





Lack of Inclusiveness

- Lack of diversity
 - Of opinions
 - Of cultures
 - Of backgrounds



Lack of Inclusiveness

- Lack of diversity
 - Of opinions
 - Of cultures
 - Of backgrounds
- Who gets to decide
 - How does state-of-the-art look like?
 - What to work on?
 - Which data to use?



Lack of Inclusiveness

- Lack of diversity
 - Of opinions
 - Of cultures
 - Of backgrounds
- Who gets to decide
 - How does state-of-the-art look like?
 - What to work on?
 - Which data to use?
- Bad science!



Daniel Khashabi @DanielKhashabi

...

It is concerning that an increasing number of research papers base the core of their studies/findings on the new GPT3 models (especially 'davinci-002'), which we know little about training/tuning. How can we do scientific research on these murky foundations?

2:46 PM · Oct 25, 2022 · Twitter Web App

14 Retweets 5 Quote Tweets 181 Likes



Lack of Inclusiveness

- Lack of diversity
 - Of opinions
 - Of cultures
 - Of backgrounds
- Who gets to decide
 - How does state-of-the-art look like?
 - What to work on?
 - Which data to use?
- Bad science!



Daniel Khashabi @DanielKhashabi

...

It is concerning that an increasing number of research papers base the core of their studies/findings on the new GPT3 models (especially 'davinci-002'), which we know little about training/tuning. How can we do scientific research on these murky foundations?

2:46 PM · Oct 25, 2022 · Twitter Web App

14 Retweets 5 Quote Tweets 181 Likes



Percy Liang @percyliang · 16h

Nice to see more info on OpenAI's models: beta.openai.com/docs/model-index...

TIL text-davinci-002 is actually based on code-davinci-002! It's also remarkable how far the models described in papers are from the APIs that everyone's using...

1

6

59

↑

Energy and Policy Considerations for Deep Learning in NLP

ACL 2019



Emma Strubell

Ananya Ganesh

Andrew McCallum

College of Information and Computer Sciences

University of Massachusetts Amherst

{strubell, aganesh, mccallum}@cs.umass.edu

Consumption	CO₂e (lbs)
Air travel, 1 person, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

Training one model (GPU)

NLP pipeline (parsing, SRL)	39
w/ tuning & experiments	78,468
Transformer (big)	192
w/ neural arch. search	626,155

Energy and Policy Considerations for Deep Learning in NLP

ACL 2019



Emma Strubell

Ananya Ganesh

Andrew McCallum

College of Information and Computer Sciences

University of Massachusetts Amherst

{strubell, aganesh, mccallum}@cs.umass.edu

Consumption	CO₂e (lbs)
Air travel, 1 person, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

Training one model (GPU)

NLP pipeline (parsing, SRL) w/ tuning & experiments	39 78,468
Transformer (big)	192
w/ neural arch. search	626,155

*Is AI really creating an
environmental problem?*



Google's Answer: No!

BLOG ›

Good News About the Carbon Footprint of Machine Learning Training

TUESDAY, FEBRUARY 15, 2022

Posted by David Patterson, Distinguished Engineer, Google Research, Brain Team

Strubell et al.'s energy estimate for NAS ended up 18.7X too high for the average organization (...) and 88X off in emissions for energy-efficient organizations like Google

Our Answer: Maybe?

Measuring the Carbon Intensity of AI in Cloud Instances

JESSE DODGE, Allen Institute for AI, USA

TAYLOR PREWITT, University of Washington, USA

REMI TACHET DES COMBES, Microsoft Research Montreal, USA

ERIKA ODMARK, Microsoft, USA

ROY SCHWARTZ, Hebrew University of Jerusalem, Israel

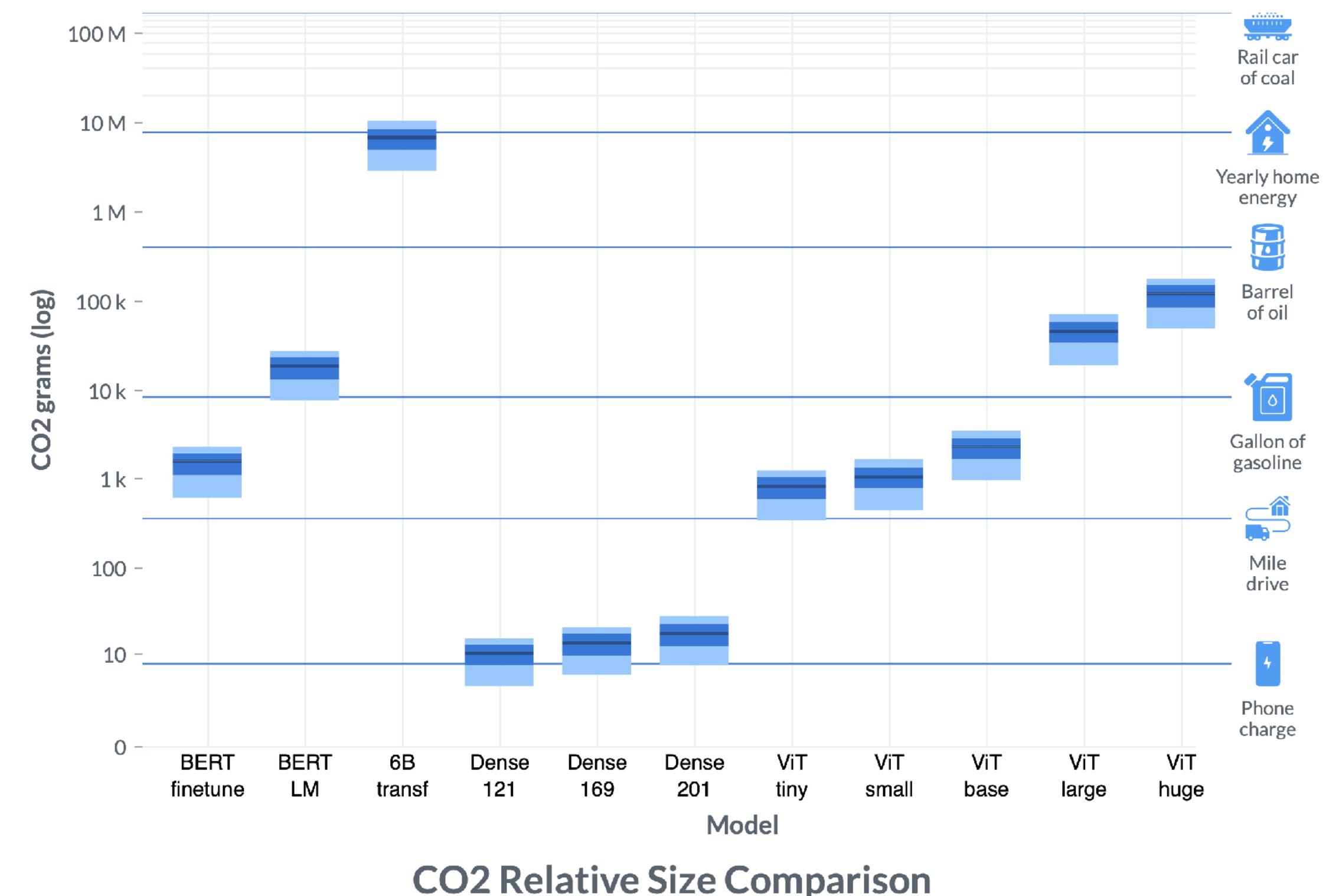
EMMA STRUBELL, Carnegie Mellon University, USA

ALEXANDRA SASHA LUCCIONI, Hugging Face, USA

NOAH A. SMITH, Allen Institute for AI and University of Washington, USA

NICOLE DECARIO, Allen Institute for AI, USA

WILL BUCHANAN, Microsoft, USA



Green
Software
Foundation





Meas

JESSE

TAYLO

REMIT

ERIKA

ROY SC

EMMA

ALEXA

NOAH

NICOL

WILL B

CO₂ grams (log)

100 M

10 M

1 M

100 k

10 k

1 k

100

10

0

BERT
finetune

BERT
LM

6B
transf

Dense
121

Dense
169

Dense
201

ViT
tiny

ViT
small

ViT
base

ViT
large

ViT
huge

Model



Rail car
of coal



Yearly home
energy



Barrel
of oil



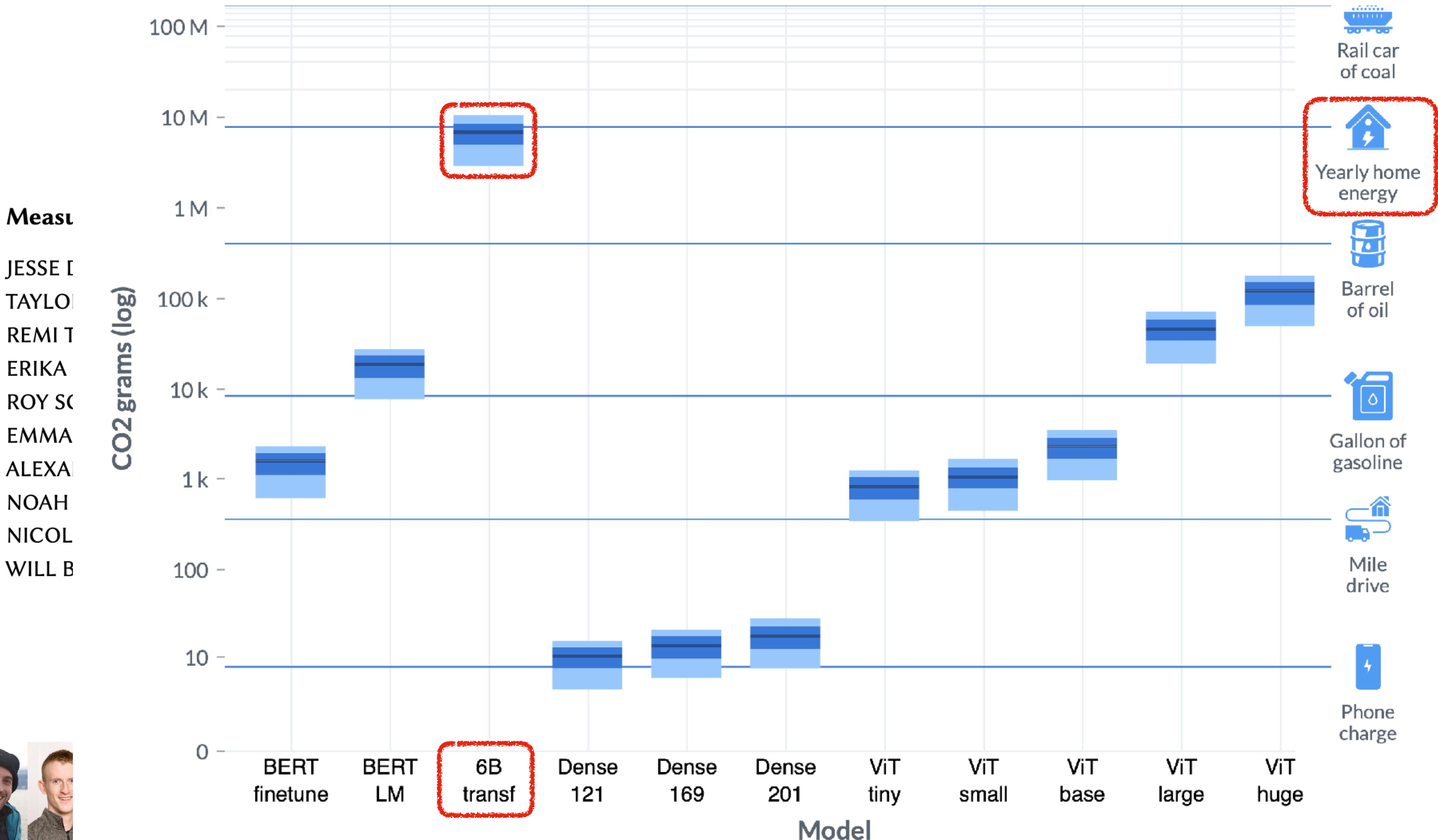
Gallon of
gasoline



Mile
drive



Phone
charge





AI and the Environment

- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - These are typically run very few times



AI and the Environment

- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - These are typically run very few times
- What about “normal” experiments?
 - Much **cheaper**, but run **hundreds / thousands of times a day?**



AI and the Environment

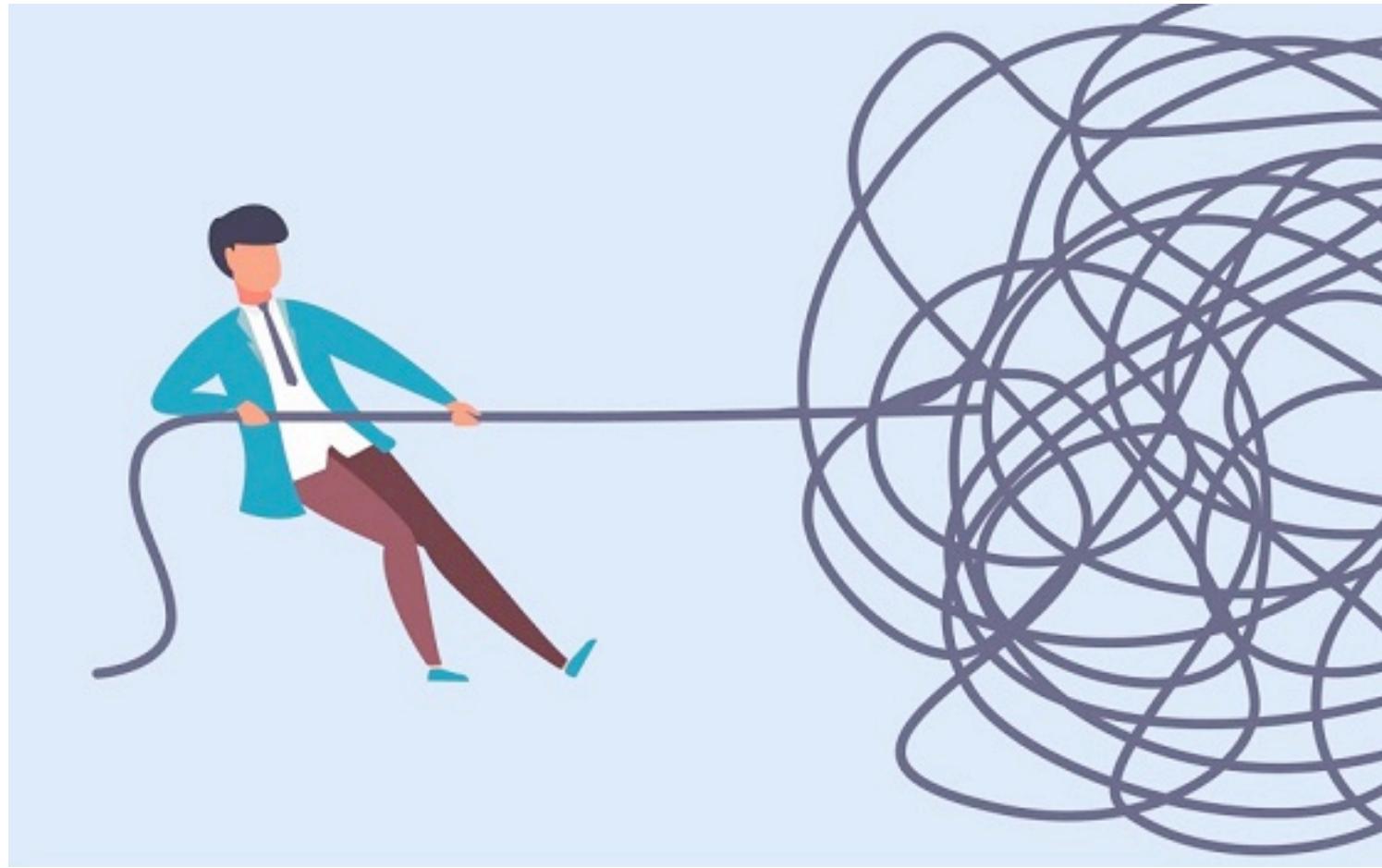
- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - These are typically run very few times
- What about “normal” experiments?
 - Much **cheaper**, but run **hundreds / thousands of times a day?**
- What about inference operations?
 - Very **cheap** (though increasingly more expensive)
 - Run **billions of times a day?**
 - 80-90% of AI computation is spent on inference



AI and the Environment

- Evidence around the **most expensive experiments**
 - More recent models consume 2-3 orders of magnitude more CO₂ (Luccioni et al., 2022)
 - These are typically run very few times
- What about “normal” experiments?
 - Much **cheaper**, but run **hundreds / thousands of times a day?**
- What about inference operations?
 - Very **cheap** (though increasingly more expensive)
 - Run **billions of times a day?**
 - 80-90% of AI computation is spent on inference
- Much harder to **measure!**

Outline



Challenges

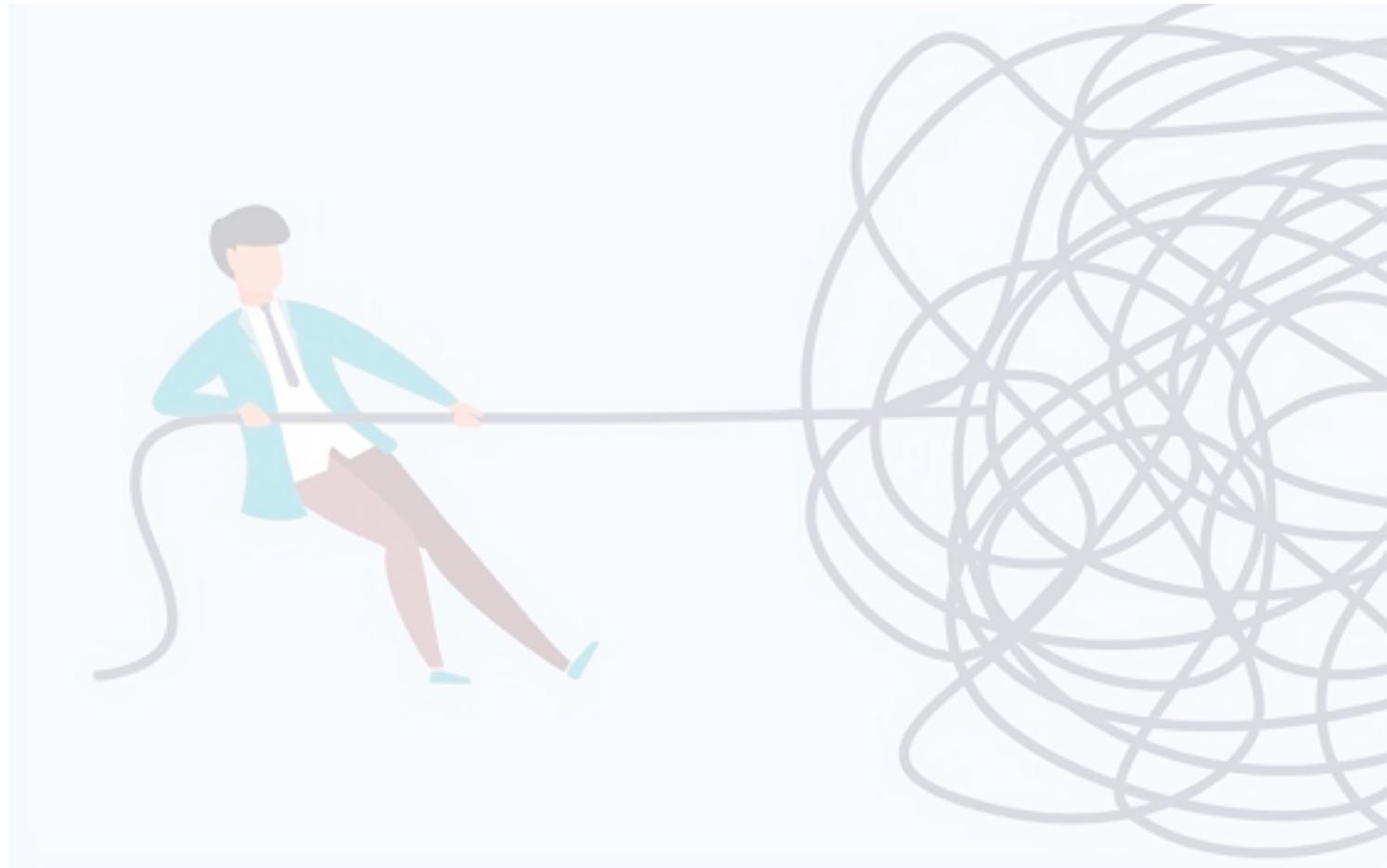


Opportunities



Mitigation

Outline



Challenges



Opportunities



Mitigation



Stop training *large* models?



Large Models are Important

- Push the limits of SOTA
- Released large pre-trained models **save compute**
- Large models are potentially faster to train
 - Li et al. (2020)



Large Models are Important

- Push the limits of SOTA
- Released large pre-trained models **save compute**
- Large models are potentially faster to train
 - Li et al. (2020)
- But, **large models have concerning side affects**
 - Inclusiveness, environment



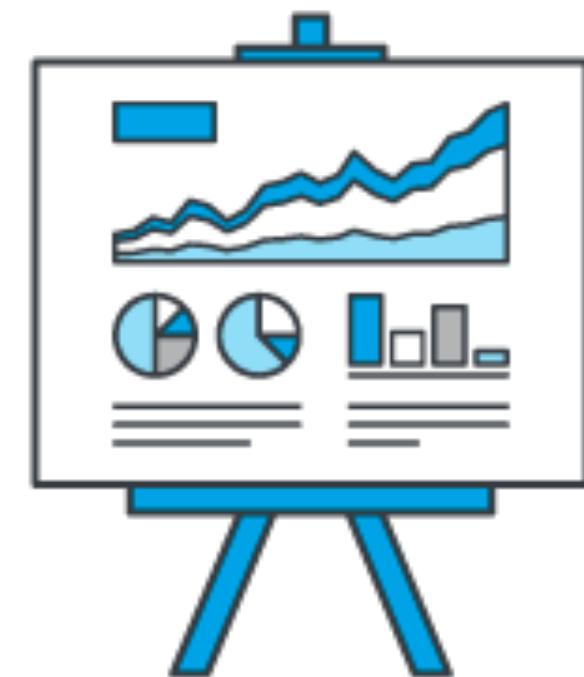
Large Models are Important

- Push the limits of SOTA
- Released large pre-trained models **save compute**
- Large models are potentially faster to train
 - Li et al. (2020)
- But, **large models have concerning side affects**
 - Inclusiveness, environment
- Our goal is to **mitigate these side affects**

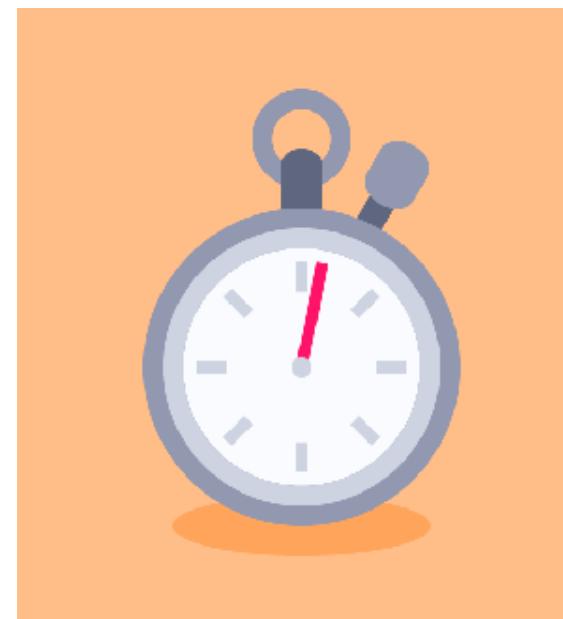
Mitigating the Challenges of Scaling



**Enhanced
Reporting**



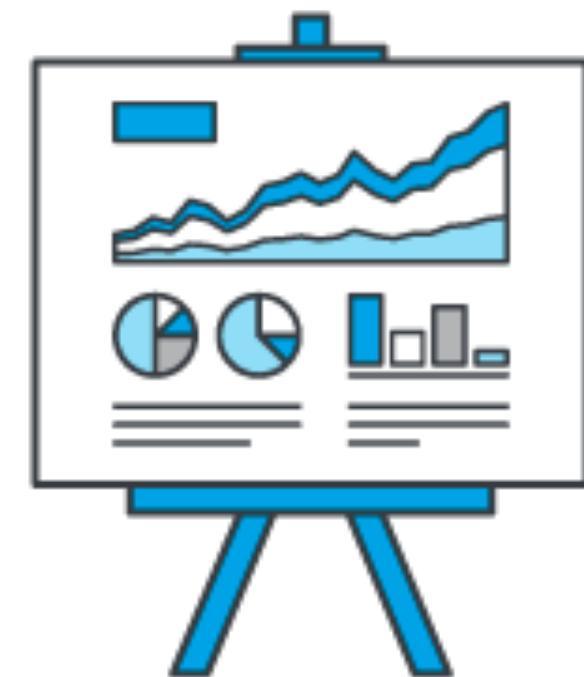
**Efficient
Methods**



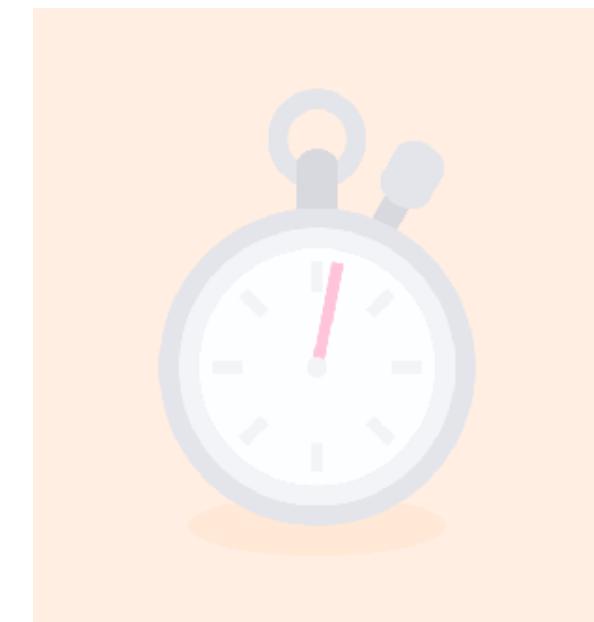
Mitigating the Challenges of Scaling



**Enhanced
Reporting**



**Efficient
Methods**





Is Model A > Model B?

Reimers & Gurevych (2017)

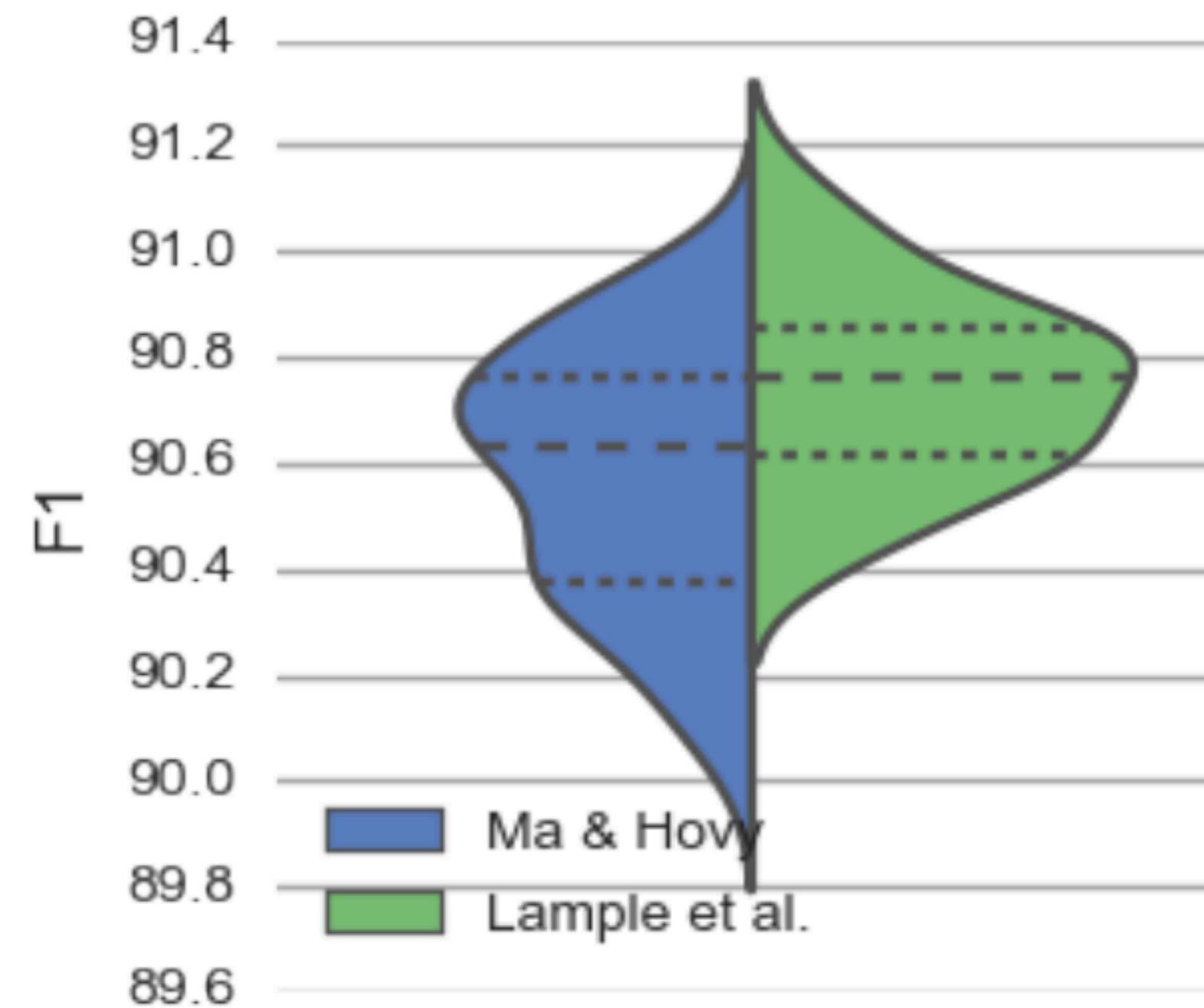
Model	F1
Model A	91.21
Model B	90.94



Is Model A > Model B?

Reimers & Gurevych (2017)

Model	F1
Model A	91.21
Model B	90.94





Is Model A > Model B?

Melis et al. (2018)

Model	Size	Depth	Valid	Test	Perplexity (↓)
Medium LSTM, Zaremba et al. (2014)	10M	2	86.2	82.7	
Large LSTM, Zaremba et al. (2014)	24M	2	82.2	78.4	
VD LSTM, Press & Wolf (2016)	51M	2	75.8	73.2	
VD LSTM, Inan et al. (2016)	9M	2	77.1	73.9	
VD LSTM, Inan et al. (2016)	28M	2	72.5	69.0	
VD RHN, Zilly et al. (2016)	24M	10	67.9	65.4	
NAS, Zoph & Le (2016)	25M	-	-	64.0	
NAS, Zoph & Le (2016)	54M	-	-	62.4	
AWD-LSTM, Merity et al. (2017) †	24M	3	60.0	57.3	



Is Model A > Model B?

Melis et al. (2018)

Model	Size	Depth	Valid	Test	Perplexity (↓)
Medium LSTM, Zaremba et al. (2014)	10M	2	86.2	82.7	
Large LSTM, Zaremba et al. (2014)	24M	2	82.2	78.4	
VD LSTM, Press & Wolf (2016)	51M	2	75.8	73.2	
VD LSTM, Inan et al. (2016)	9M	2	77.1	73.9	
VD LSTM, Inan et al. (2016)	28M	2	72.5	69.0	
VD RHN, Zilly et al. (2016)	24M	10	67.9	65.4	
NAS, Zoph & Le (2016)	25M	-	-	64.0	
NAS, Zoph & Le (2016)	54M	-	-	62.4	
AWD-LSTM, Merity et al. (2017) †	24M	3	60.0	57.3	



Is Model A > Model B?

Melis et al. (2018)

Model	Size	Depth	Valid	Test	Perplexity (↓)
Medium LSTM, Zaremba et al. (2014)	10M	2	86.2	82.7	
Large LSTM, Zaremba et al. (2014)	24M	2	82.2	78.4	
VD LSTM, Press & Wolf (2016)	51M	2	75.8	73.2	
VD LSTM, Inan et al. (2016)	9M	2	77.1	73.9	
VD LSTM, Inan et al. (2016)	28M	2	72.5	69.0	
VD RHN, Zilly et al. (2016)	24M	10	67.9	65.4	
NAS, Zoph & Le (2016)	25M	-	-	64.0	
NAS, Zoph & Le (2016)	54M	-	-	62.4	
AWD-LSTM, Merity et al. (2017) †	24M	3	60.0	57.3	
<hr/>					
LSTM		1	61.8	59.6	
LSTM		2	63.0	60.8	
LSTM	10M	4	62.4	60.1	
RHN		5	66.0	63.5	
NAS		1	65.6	62.7	
<hr/>					
LSTM		1	61.4	59.5	
LSTM		2	62.1	59.6	
LSTM	24M	4	60.9	58.3	
RHN		5	64.8	62.2	
NAS		1	62.1	59.7	

Carefully Tuned
(1500 trials)



Unfair Comparison

Is Model A > Model B?



Better(?) Comparison

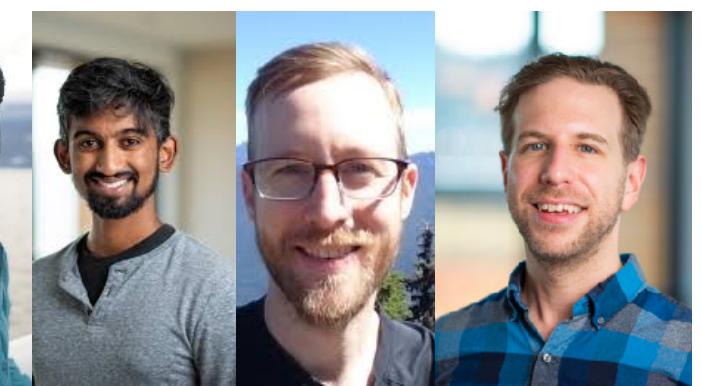
Is Model A > Model B? | *Budget*



Expected Validation

Dodge, Gururangan, Card, S. & Smith, EMNLP 2019

- Input: a set of experimental results $\{V_1, \dots, V_n\}$
- Define $V_k^* = \max_{i \in \{1, \dots, k\}} V_i$

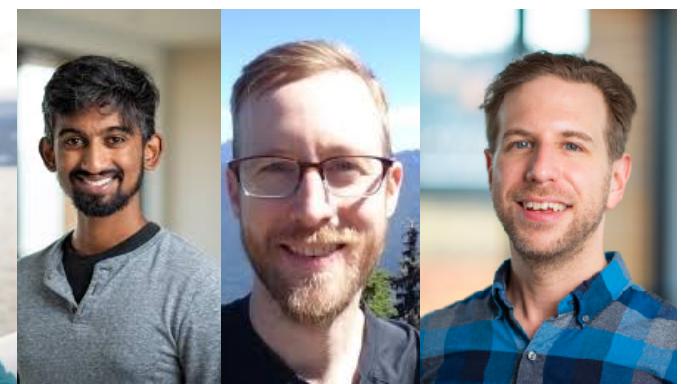




Expected Validation

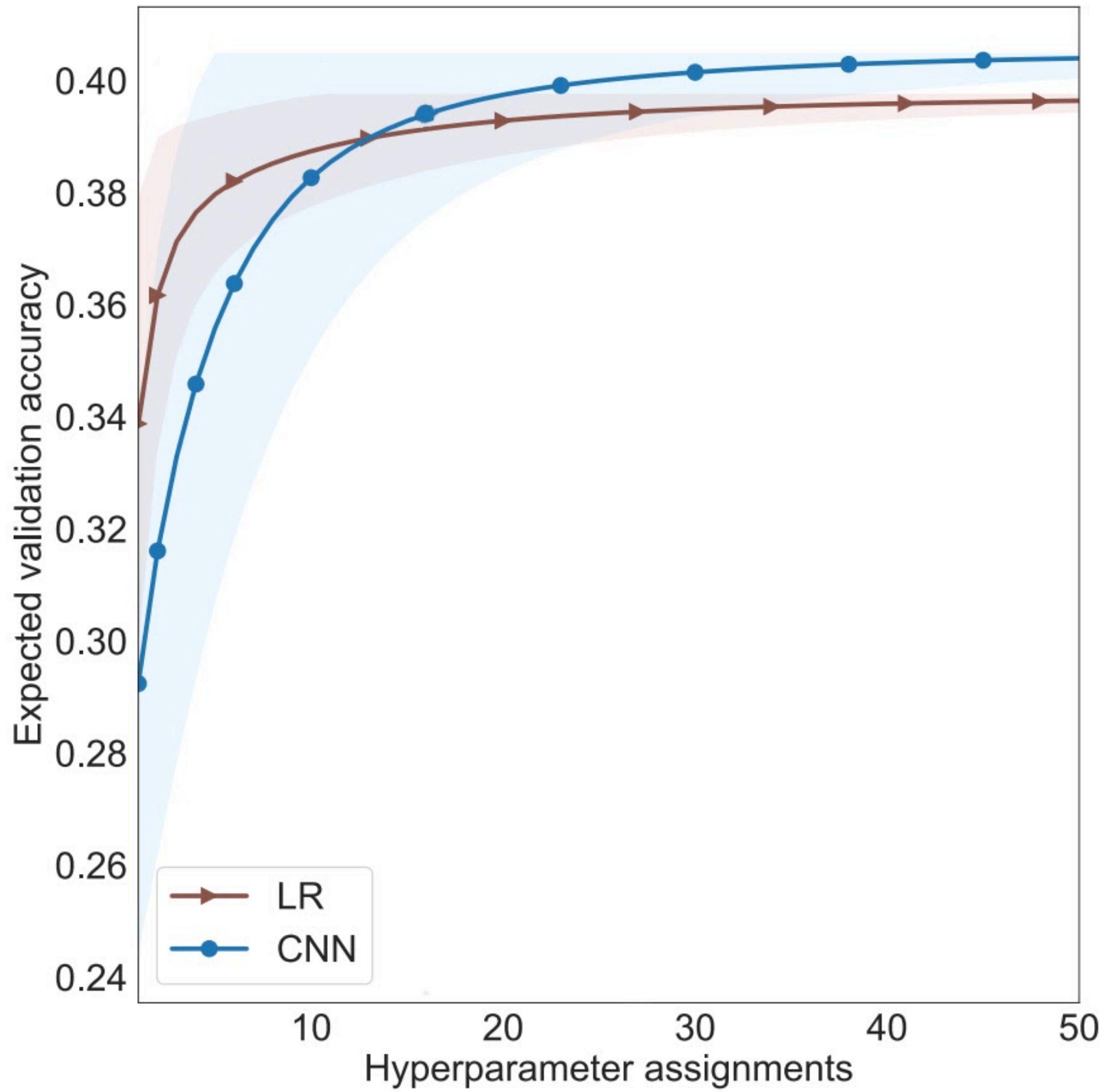
Dodge, Gururangan, Card, S. & Smith, EMNLP 2019

- Input: a set of experimental results $\{V_1, \dots, V_n\}$
- Define $V_k^* = \max_{i \in \{1, \dots, k\}} V_i$
- **Expected validation performance:** $\mathbb{E}[V_k^* | k]$
 - k=1: $mean(\{V_1, \dots, V_n\})$
 - k=2: $mean(\{\max(V_i, V_j) \forall 1 \leq i < j \leq n\})$
 - k=n: $V_n^* = \max_{i \in \{1, \dots, n\}} V_i$





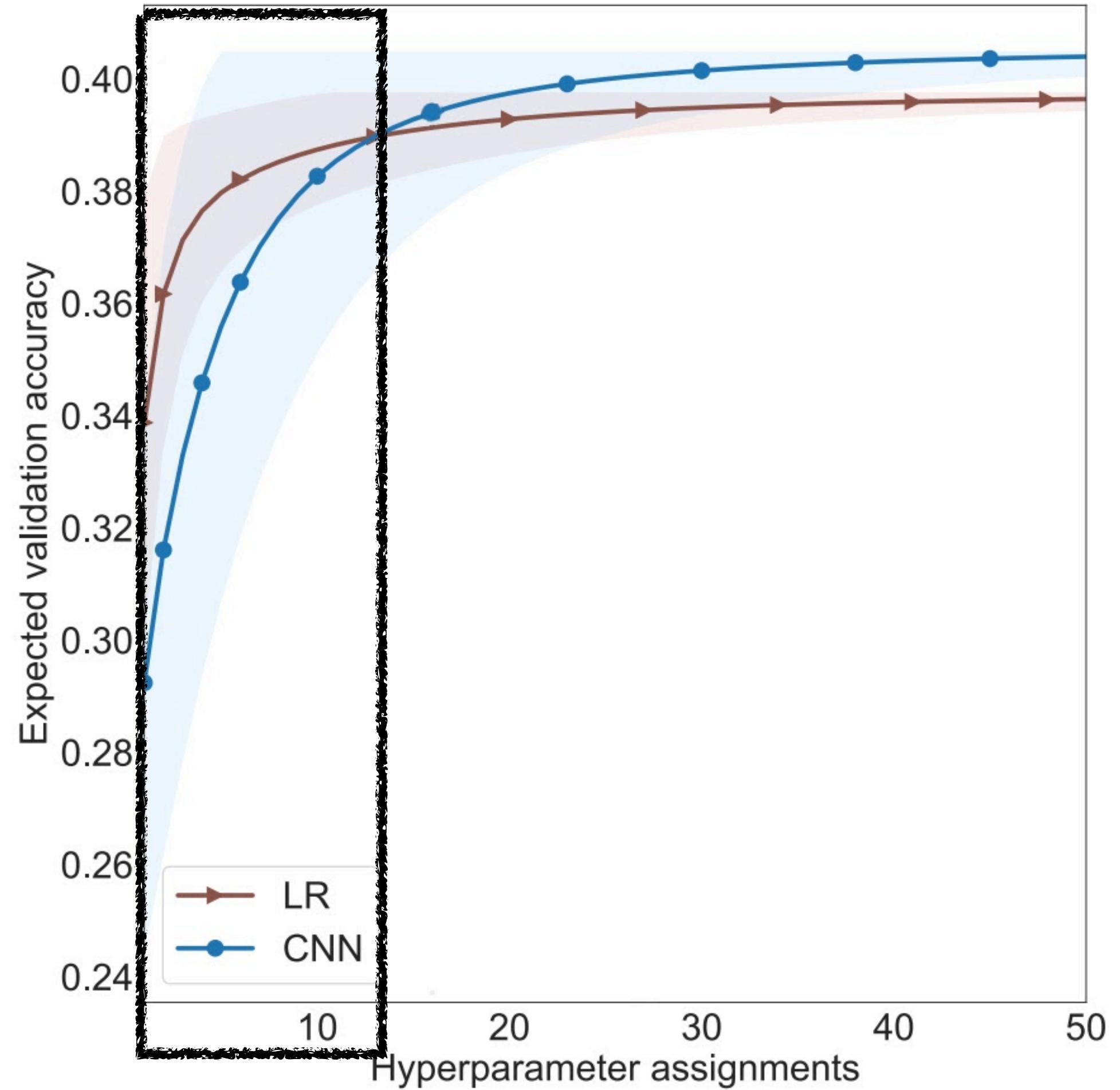
Expected Validation Example





Expected Validation Example

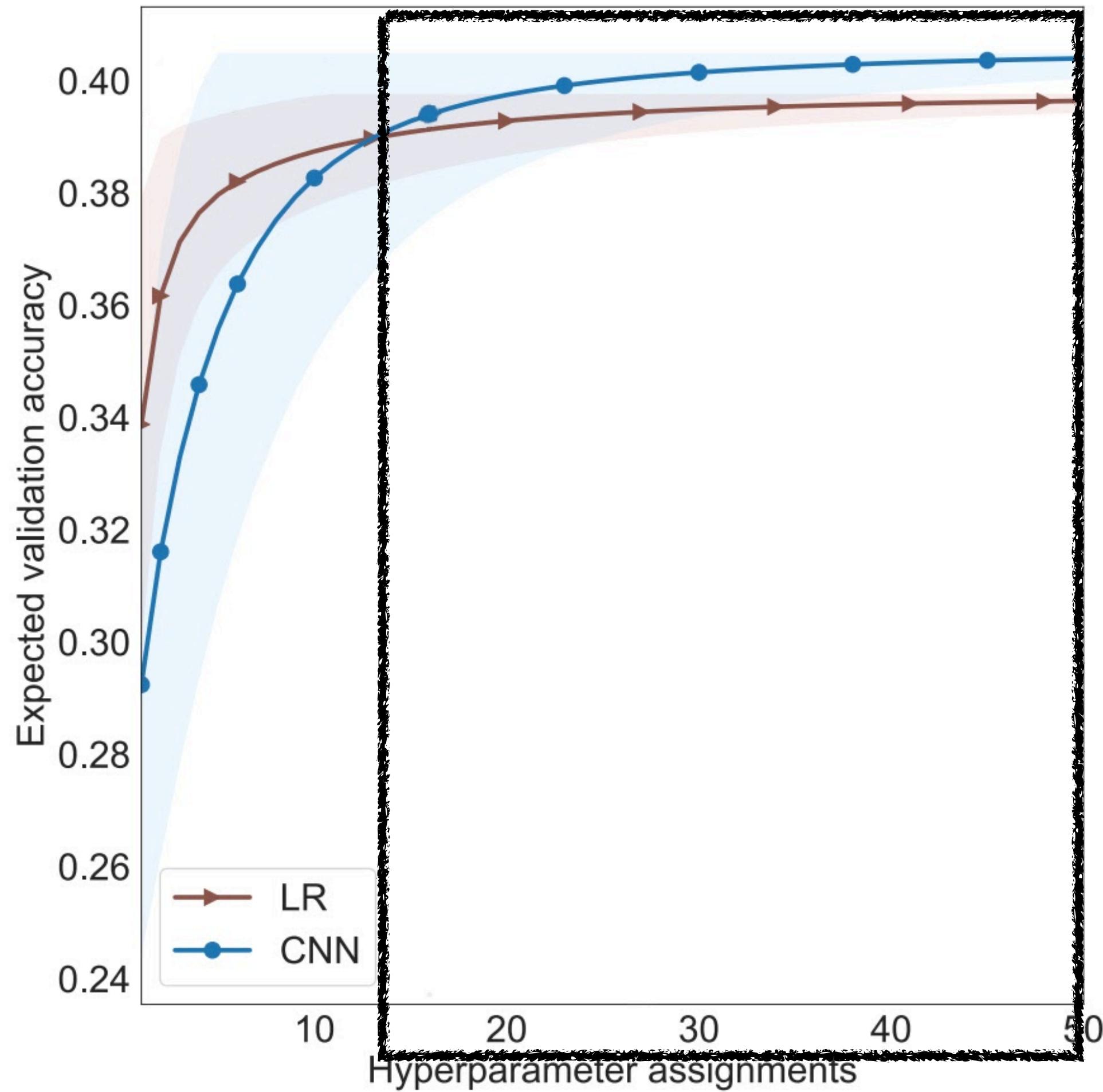
Low budget:
LR > CNN



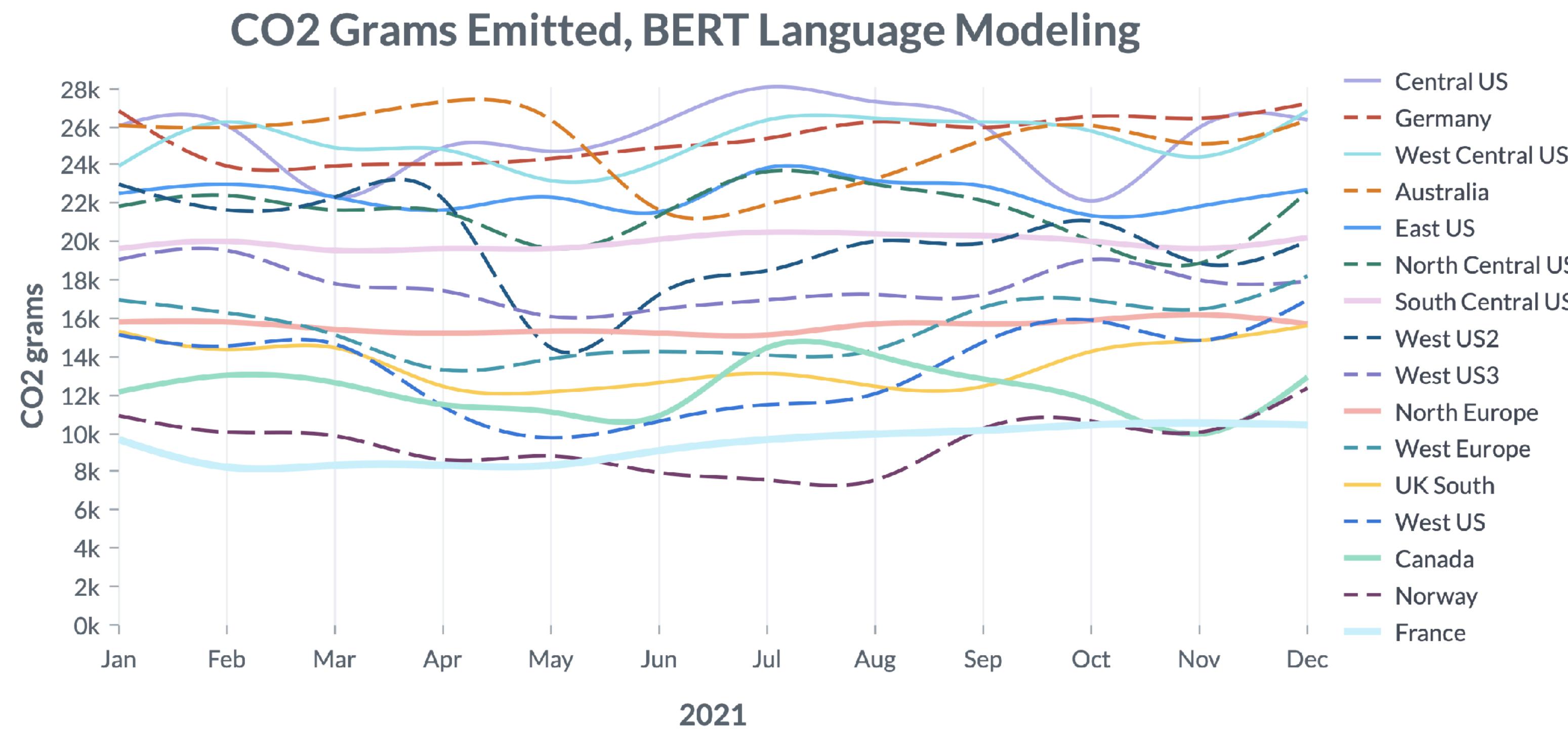


Expected Validation Example

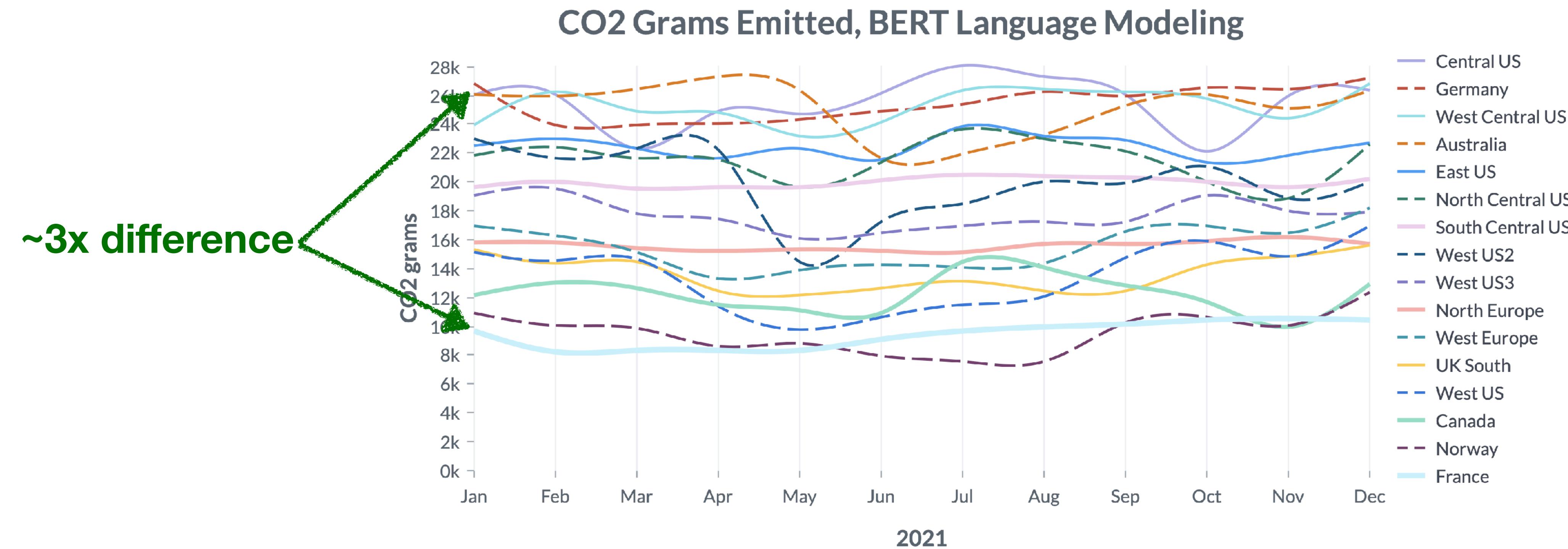
High budget:
CNN > LR



Cloud Location Matters



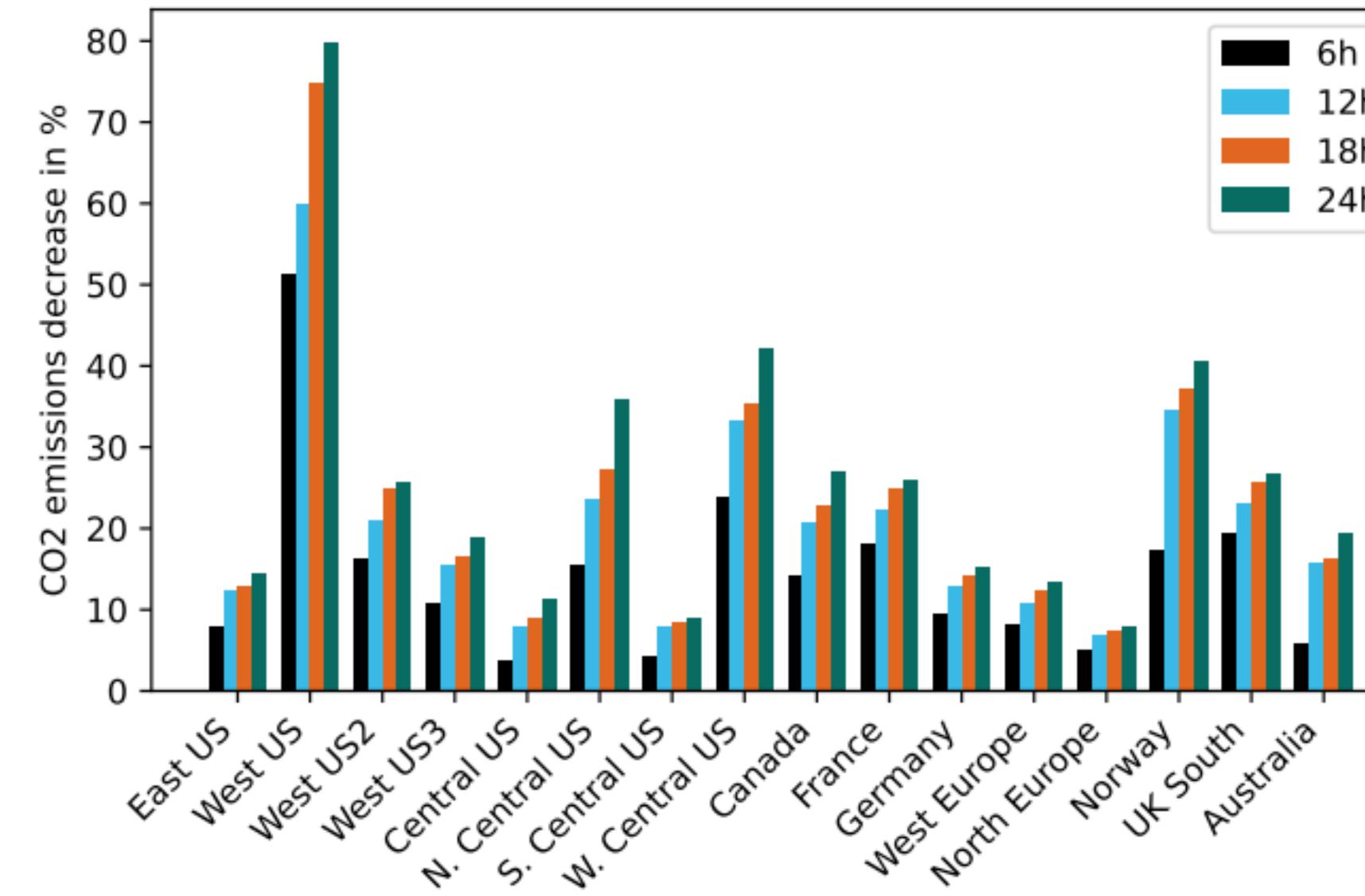
Cloud Location Matters



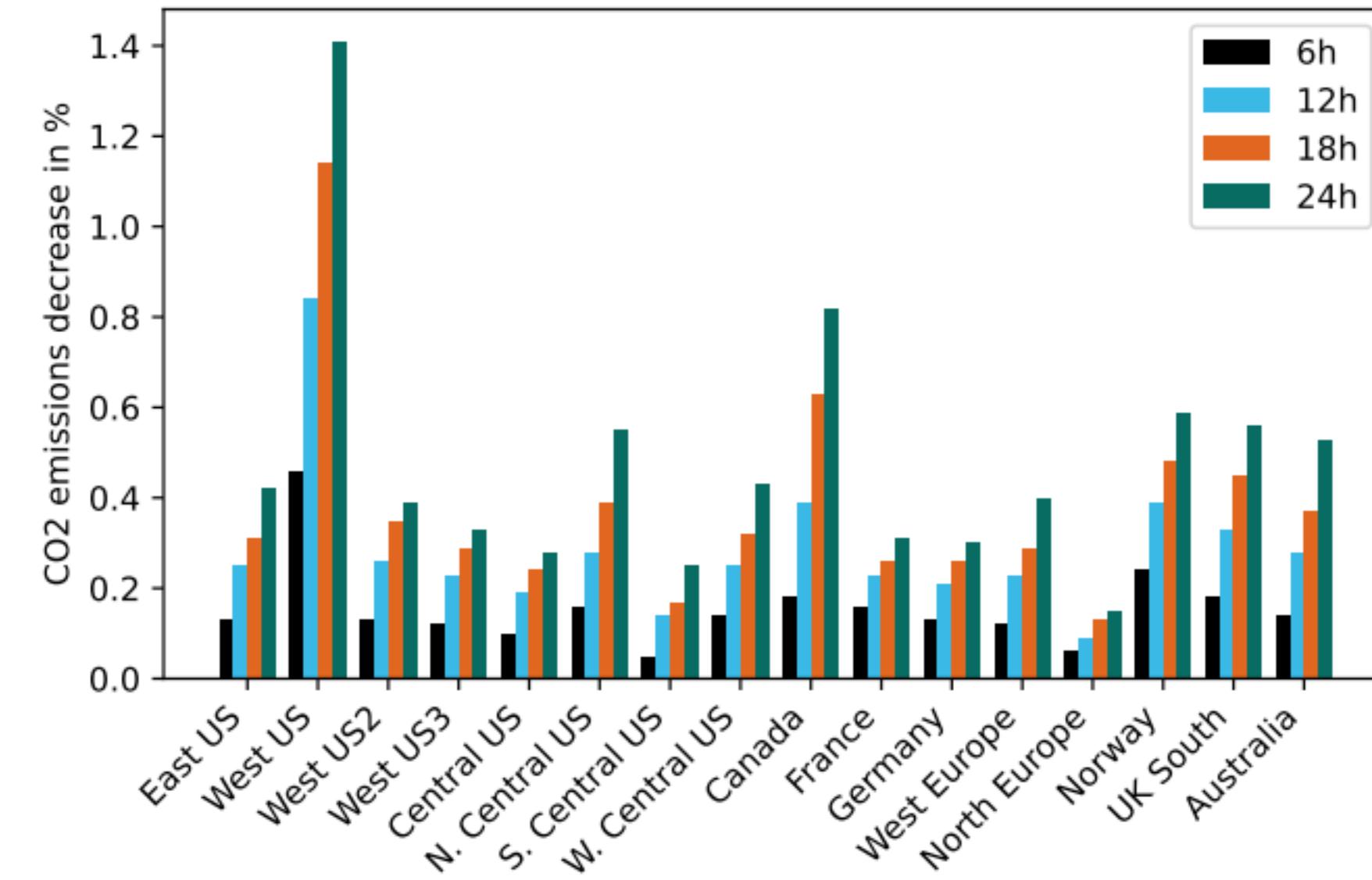


Time of Day Matters

Potential Saving with *Flexible Start*



(a) *Flexible Start* optimization for Dense 201.

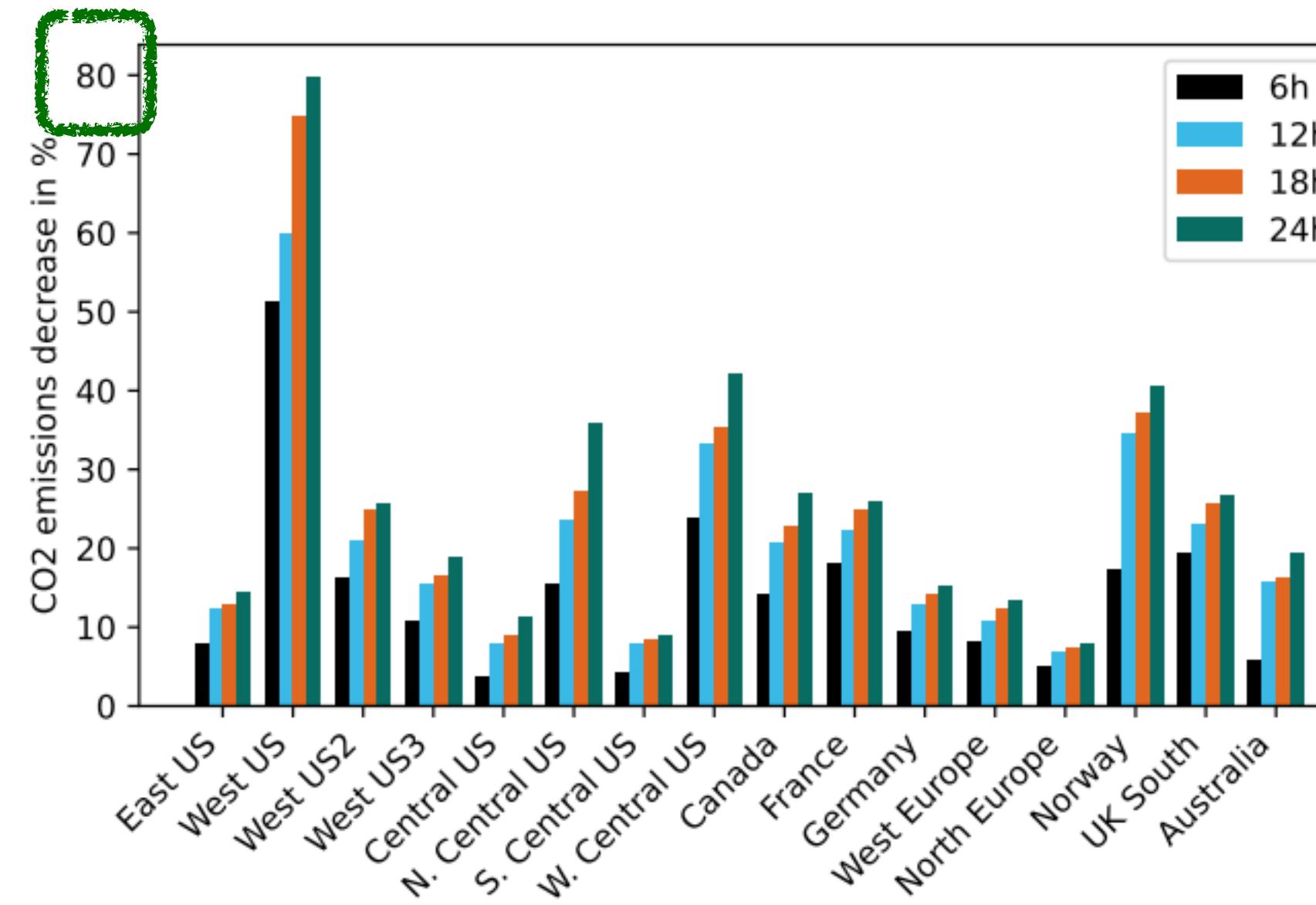


(b) *Flexible Start* optimization for 6B parameters Transformer.

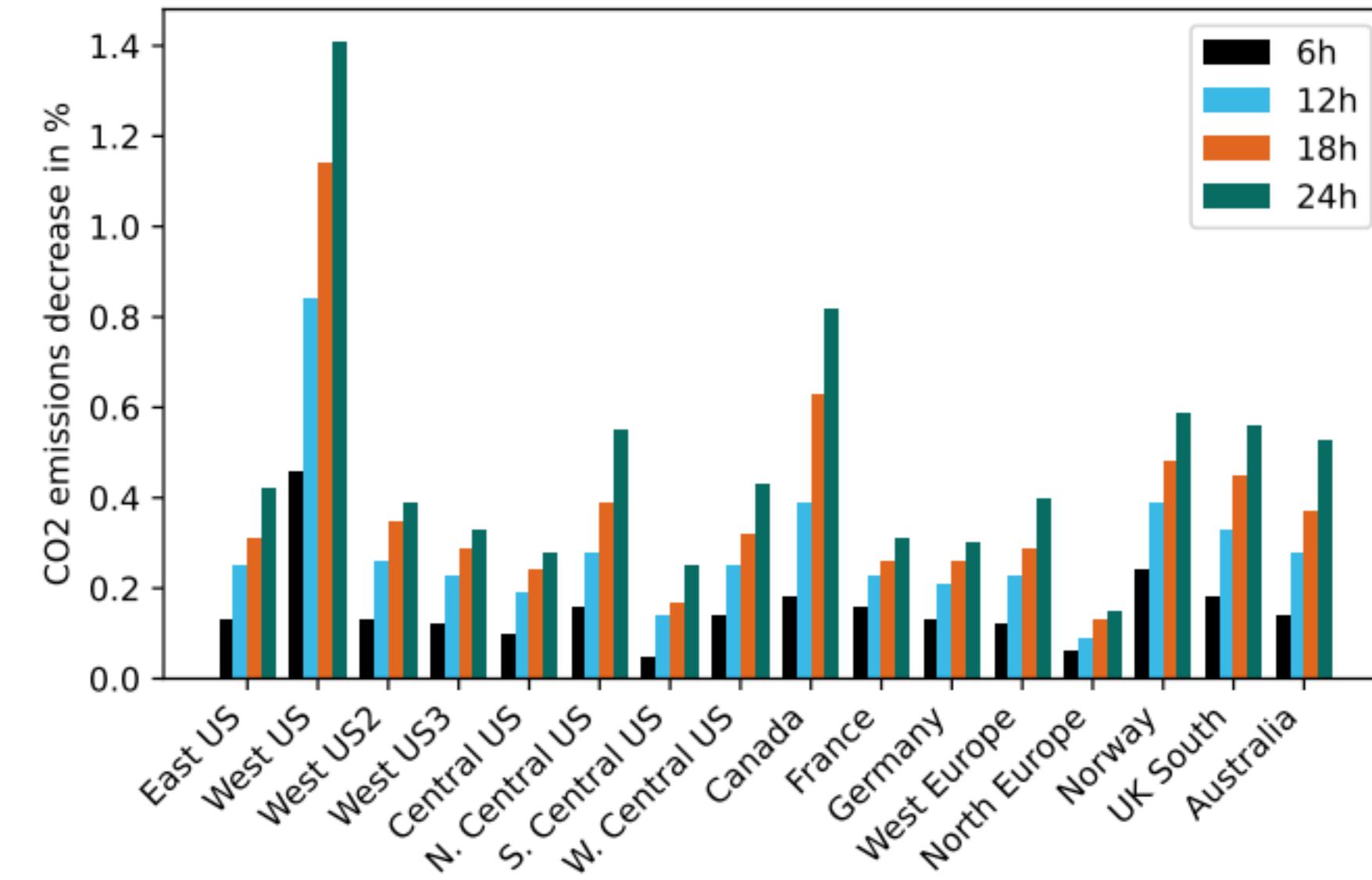


Time of Day Matters

Potential Saving with *Flexible Start*



(a) *Flexible Start* optimization for Dense 201.



(b) *Flexible Start* optimization for 6B parameters Transformer.



Reporting

Open Questions

- How much will we gain by pouring **more compute?**



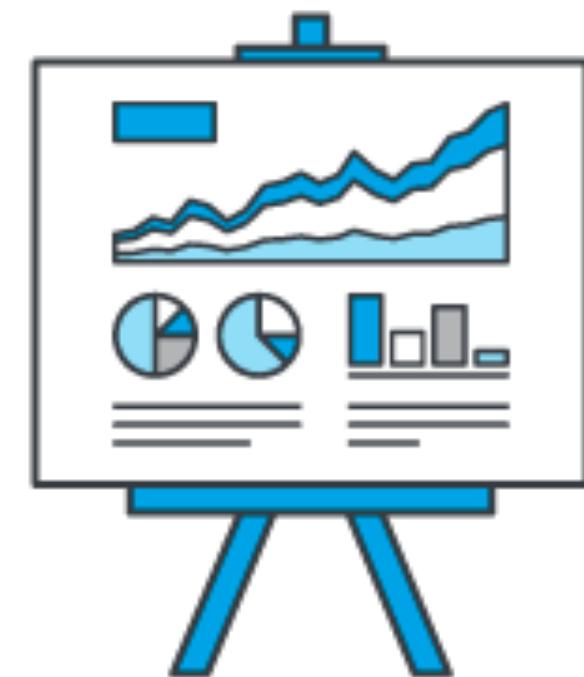
Reporting Open Questions

- How much will we gain by pouring **more compute?**
- What should we report?
 - Number of experiments
 - Time
 - FLOPs
 - Energy (KW)
 - Carbon?

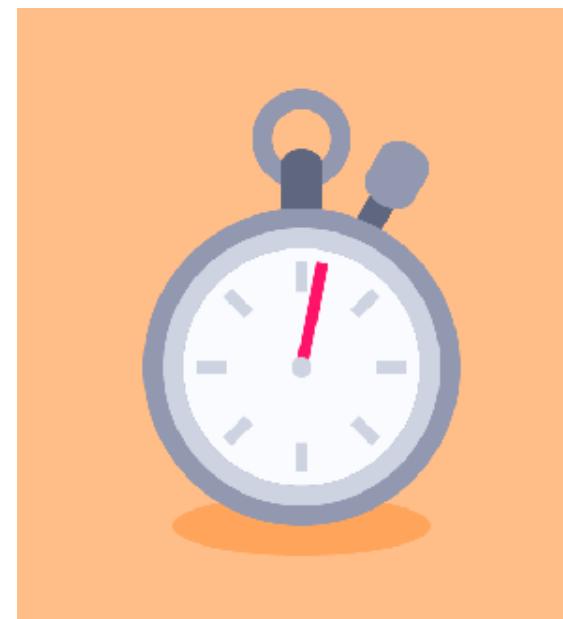
Mitigating the Challenges of Scaling



**Enhanced
Reporting**



**Efficient
Methods**



Mitigating the Challenges of Scaling



Enhanced
Reporting

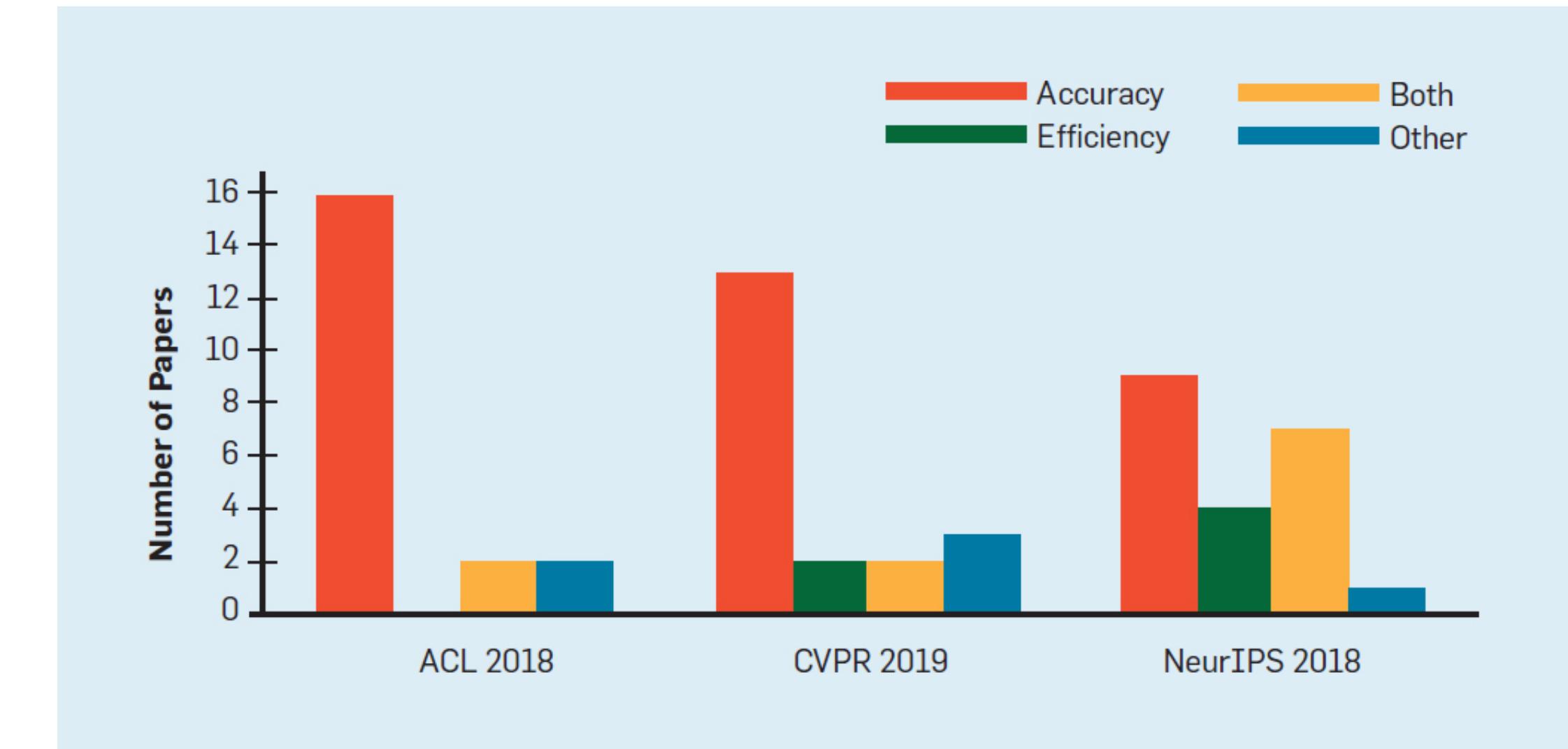


Efficient
Methods





Accuracy or Efficiency?



S. et al. (2020)



Efficient Methods for Natural Language Processing: A Survey

**Marcos Treviso^{10*}, Tianchu Ji^{3*}, Ji-Ung Lee^{7*}, Betty van Aken⁸, Qingqing Cao²,
Manuel R. Ciosici⁹, Michael Hassid¹, Kenneth Heafield¹³, Sara Hooker⁵,
Pedro H. Martins¹⁰, André F. T. Martins¹⁰, Peter Milder³, Colin Raffel⁶,**

Edwin Simpson⁴, Noam Slonim¹², Niranjan Balasubramanian³, Leon Derczynski¹¹, Roy Schwartz¹

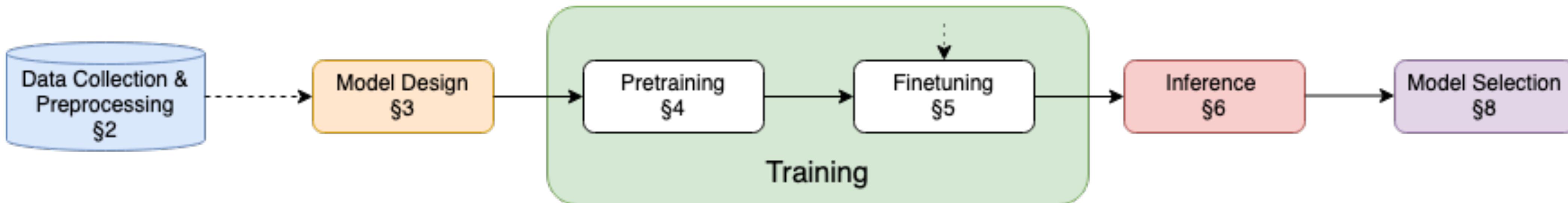
¹The Hebrew University of Jerusalem, ²University of Washington, ³Stony Brook University,

⁴University of Bristol, ⁵Cohere For AI, ⁶University of North Carolina at Chapel Hill,

⁷Technical University of Darmstadt, ⁸Berliner Hochschule für Technik,

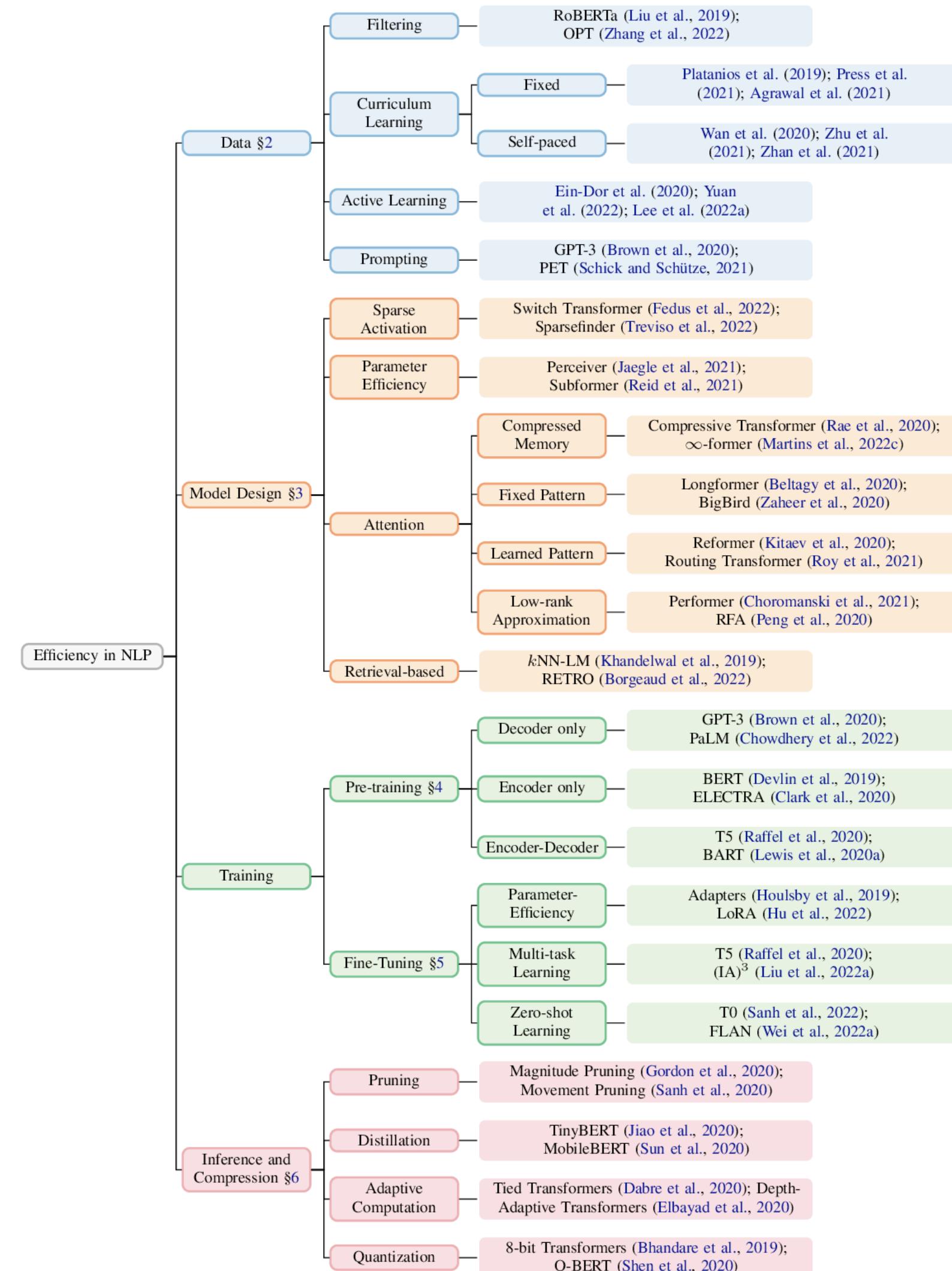
⁹University of Southern California, ¹⁰IST/University of Lisbon & Instituto de Telecomunicações,

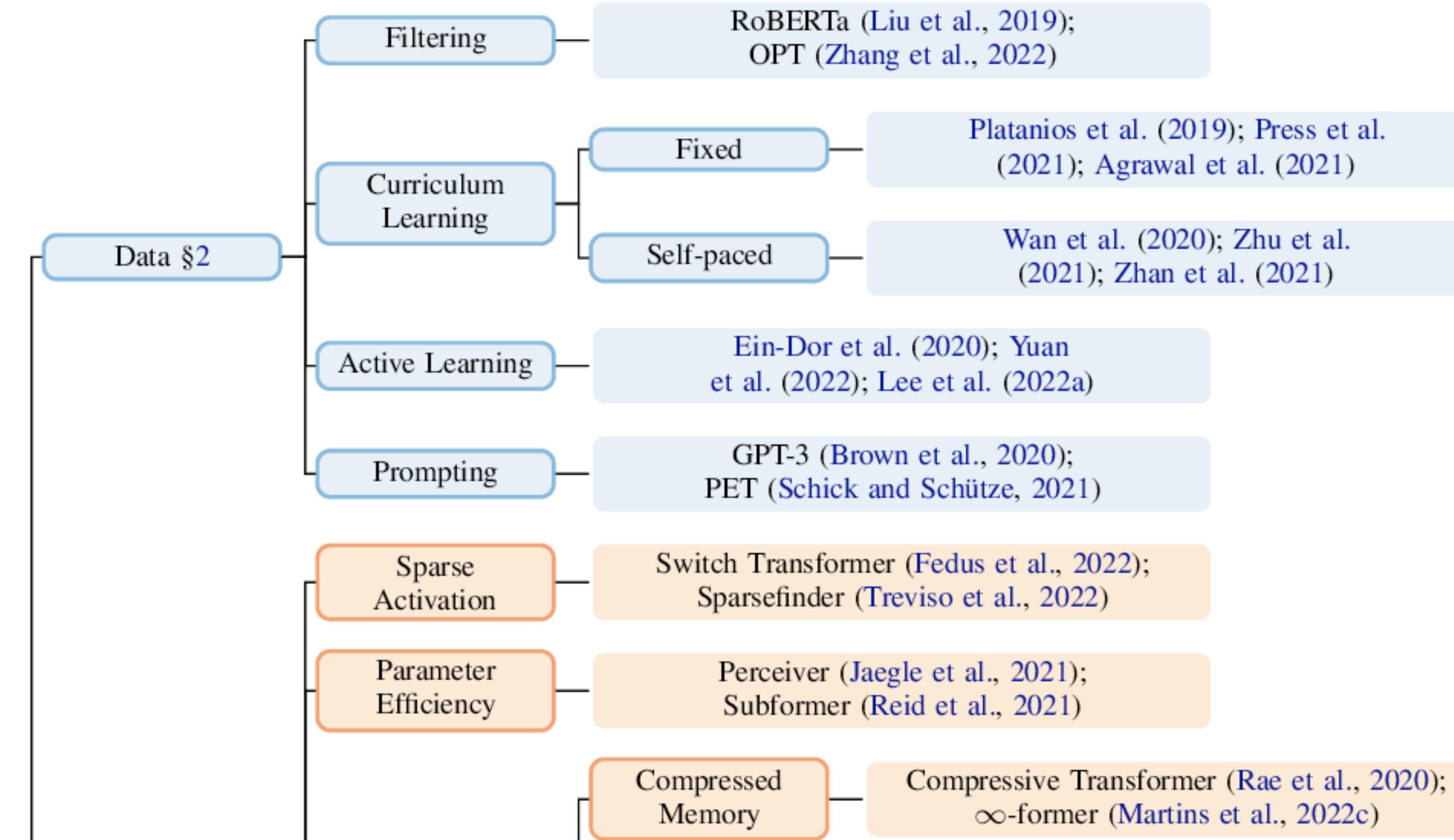
¹¹IT University of Copenhagen, ¹²IBM Research, ¹³University of Edinburgh





Efficient Methods in NLP





Filtering

- Non-text
 - Gibrish, HTML
- Text in other languages
- Foul text
- Typically done via simple, rule-based heuristics
 - Noisy process



Data-Efficient Training

Swayamdipta, S. et al., EMNLP 2020

- Not all training instances contribute the same to learning
 - Some are “easy-to-learn”, others are more challenging





Training Dynamics

- Training of neural networks typically requires iterating a large set of training examples



Training Dynamics

- Training of neural networks typically requires iterating a large set of training examples
- Models typically perform multiple iterations (aka epochs) over their training sets
 - In each of these iterations, the (partially-trained) model makes a prediction over each training instance



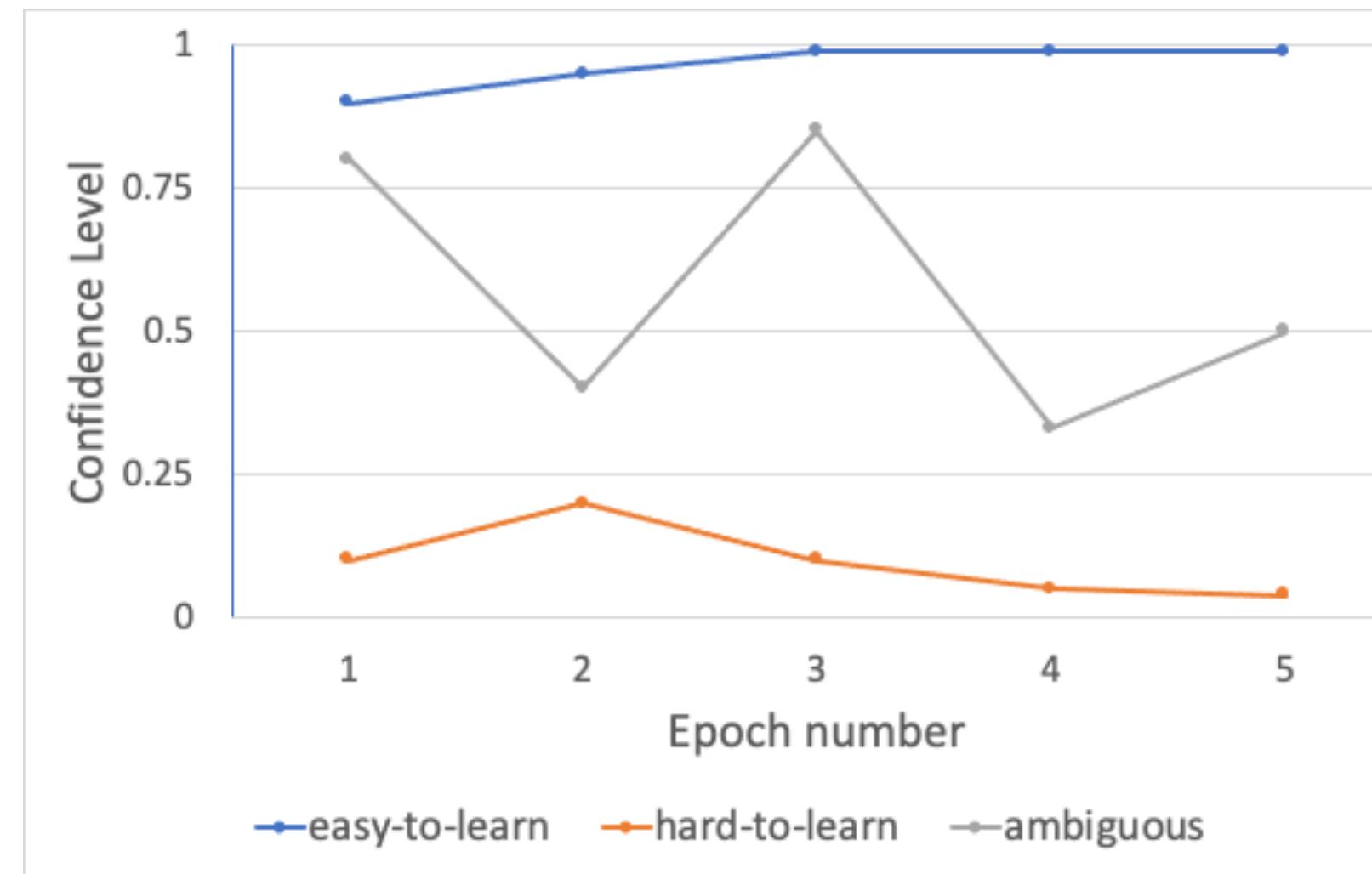
Training Dynamics

- Training of neural networks typically requires iterating a large set of training examples
- Models typically perform multiple iterations (aka epochs) over their training sets
 - In each of these iterations, the (partially-trained) model makes a prediction over each training instance
 - These predictions provide a *valuable signal*



Training Dynamics

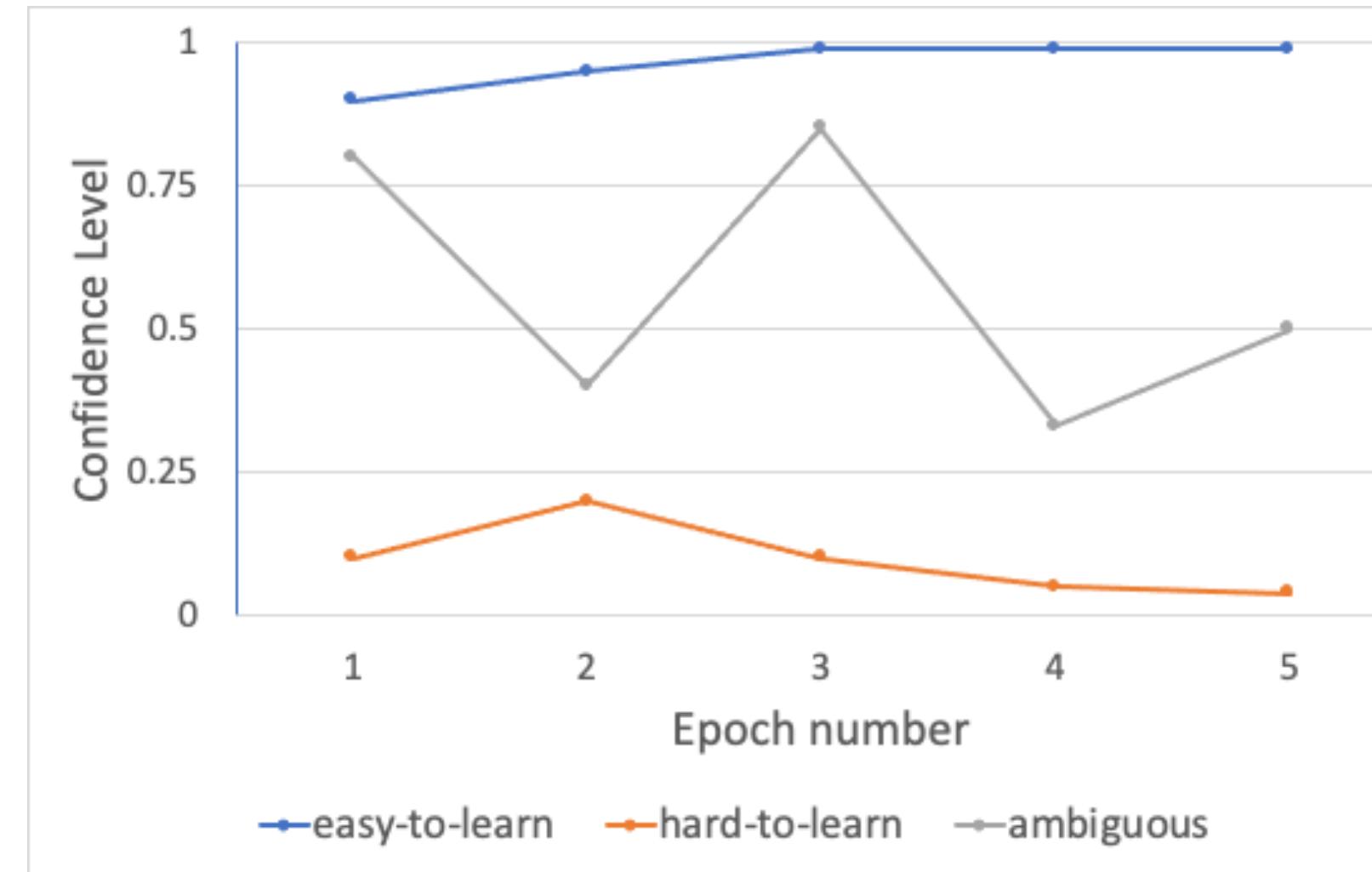
Toy Example





Training Dynamics

Toy Example

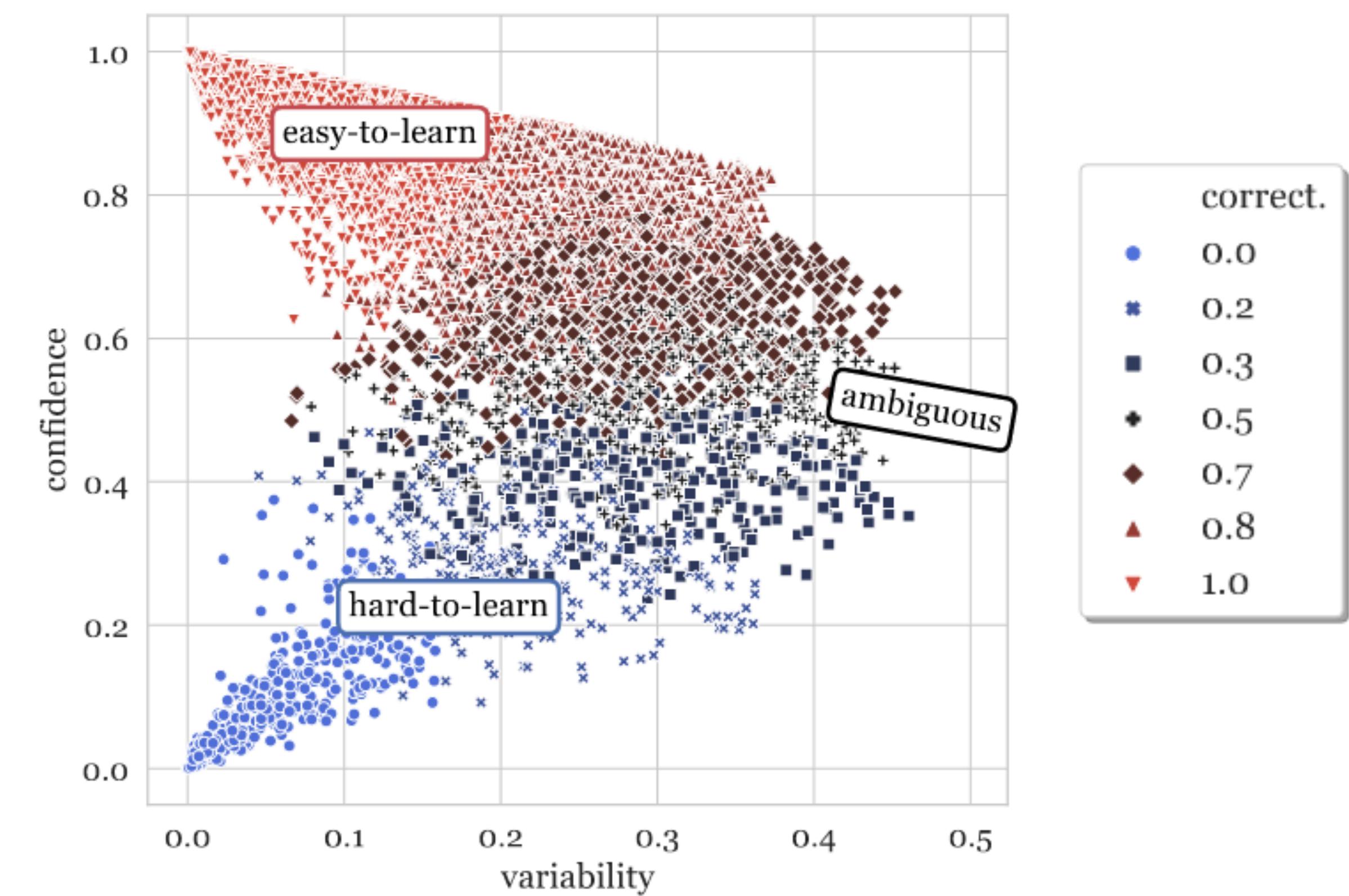


We extract two measures

- Mean confidence
- Variability

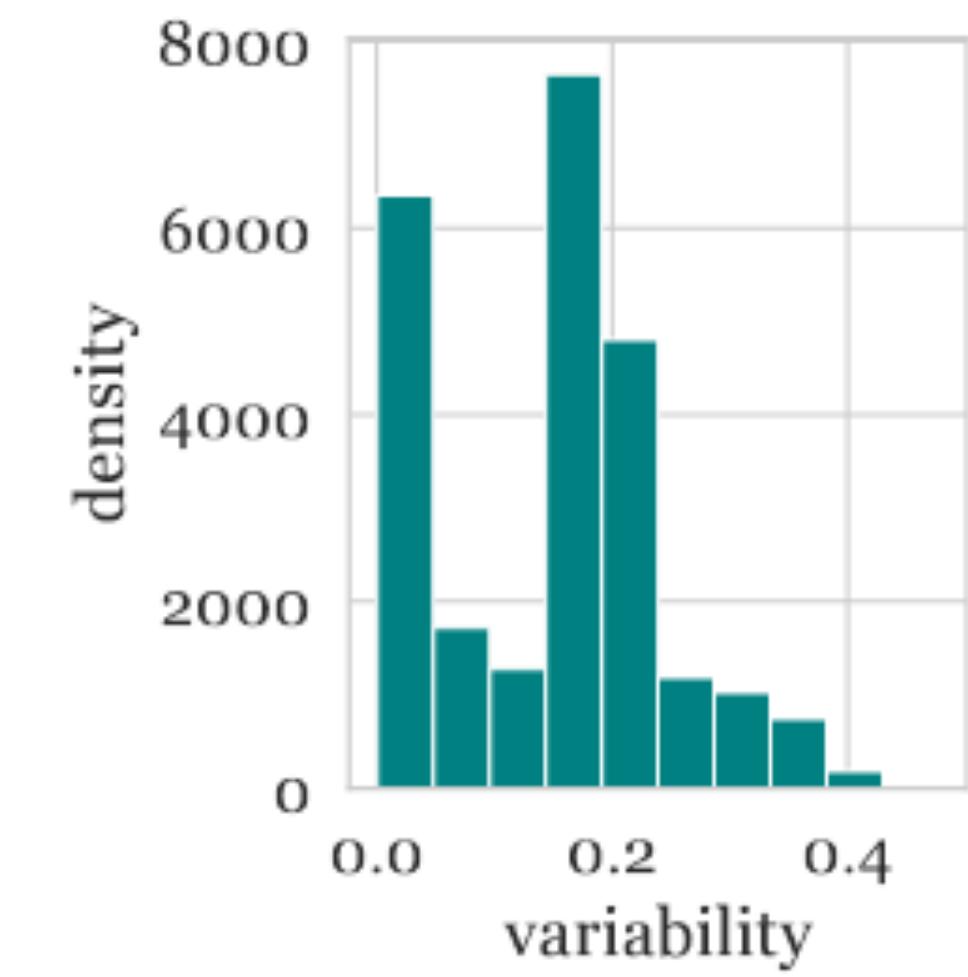
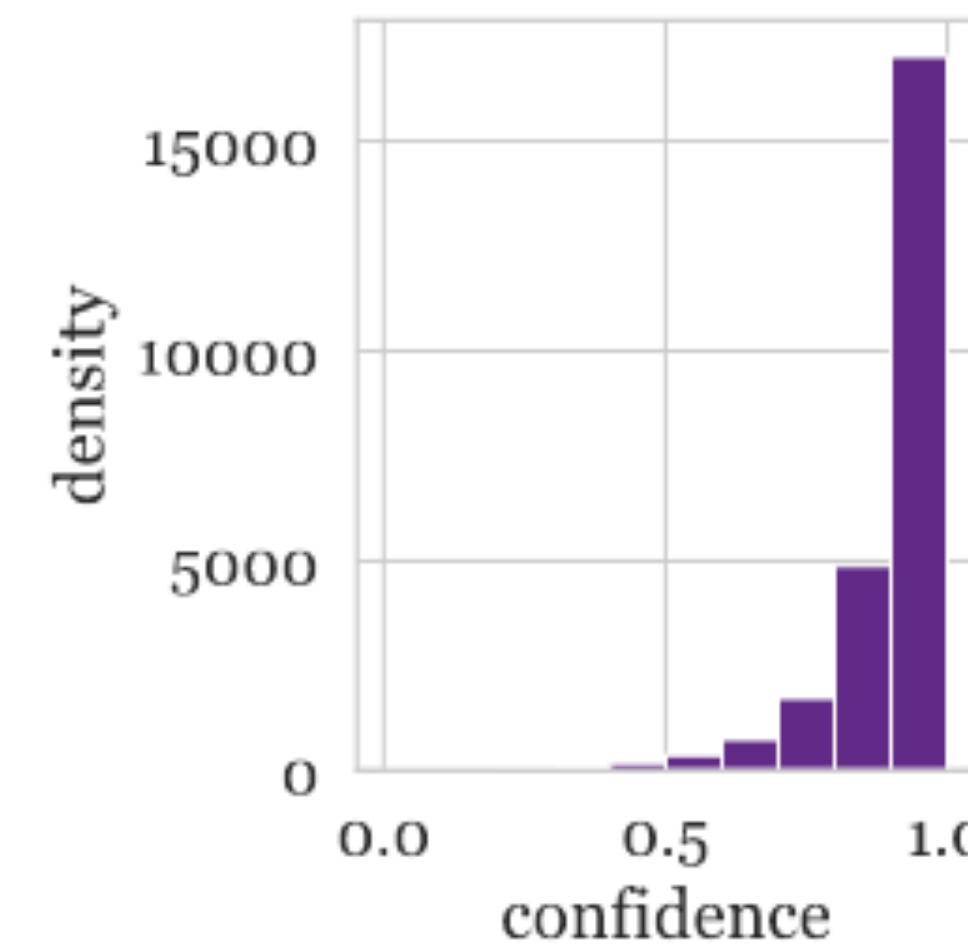


Dataset Map





Most Instances are Easy-to-Learn





Examples

	Instance	Option1	Option2
easy-to-learn	The man chose to buy the roses instead of the carnations because the __ were more beautiful.	roses*	carnations
	We enjoyed the meeting tonight but not the play as the __ was rather dull.	meeting	play*
ambiguous	The dog ran up to Leslie and away from Lawrence because __ had soap for the dog to take a bath.	Leslie ⁻	Lawrence*
	Kayla dated many more people at once than Betty, because __ was in an exclusive relationship.	Kayla*	Betty ⁺



Examples

	Instance	Option1	Option2
easy-to-learn	The man chose to buy the roses instead of the carnations because the __ were more beautiful.	roses*	carnations
	We enjoyed the meeting tonight but not the play as the __ was rather dull.	meeting	play*
ambiguous	The dog ran up to Leslie and away from Lawrence because __ had soap for the dog to take a bath.	Leslie ⁻	Lawrence*
	Kayla dated many more people at once than Betty, because __ was in an exclusive relationship.	Kayla*	Betty ⁺



Examples

	Instance	Option1	Option2
easy-to-learn	The man chose to buy the roses instead of the carnations because the __ were more beautiful.	roses*	carnations
	We enjoyed the meeting tonight but not the play as the __ was rather dull.	meeting	play*
ambiguous	The dog ran up to Leslie and away from Lawrence because __ had soap for the dog to take a bath.	Leslie ⁻	Lawrence*
	Kayla dated many more people at once than Betty, because __ was in an exclusive relationship.	Kayla*	Betty ⁺



Sample-Efficient Training

- *easy-to-learn* instances provide little value to training
- Can we use training dynamics to select the *most valuable* instances?

Experiments

WinoGrande, RoBERTa-Large

	WINOG. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

	WINOG. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

33% {

	WINOG. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

33% {

	WINOG. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

33% {

	WINOG. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

Experiments

WinoGrande, RoBERTa-Large

	WINOG. Val. (ID)	WSC (OOD)
100% train	79.7 _{0.2}	86.0 _{0.1}
random	73.3 _{1.3}	85.6 _{0.4}
<i>ambiguous</i>	78.7 _{0.4}	87.6 _{0.6}

33% {



Data Efficient Masked Language Modeling for Vision and Language

Bitton, Stanovsky, Elhadad & S., Findings of EMNLP 2021

- Given an image and a caption describing it
 - Mask some of the words in the caption
 - Train model to complete the masked words

A tiger [MASK] eating the carrot





Data Efficient Masked Language Modeling for Vision and Language

Bitton, Stanovsky, Elhadad & S., Findings of EMNLP 2021

- Given an image and a caption describing it
 - Mask some of the words in the caption
 - Train model to complete the masked words
 - Words are randomly masked with a 15% probability

A tiger [MASK] eating the carrot





Problems with Pre-training Vision and Language Models

- Of the masked tokens, roughly one half are *stop-words* or *punctuation*



Problems with Pre-training Vision and Language Models

- Of the masked tokens, roughly one half are *stop-words* or *punctuation*
- Image captions are typically short (6-7 words)
 - ⇒ in roughly 1/3 of the captions, **no word is masked!**

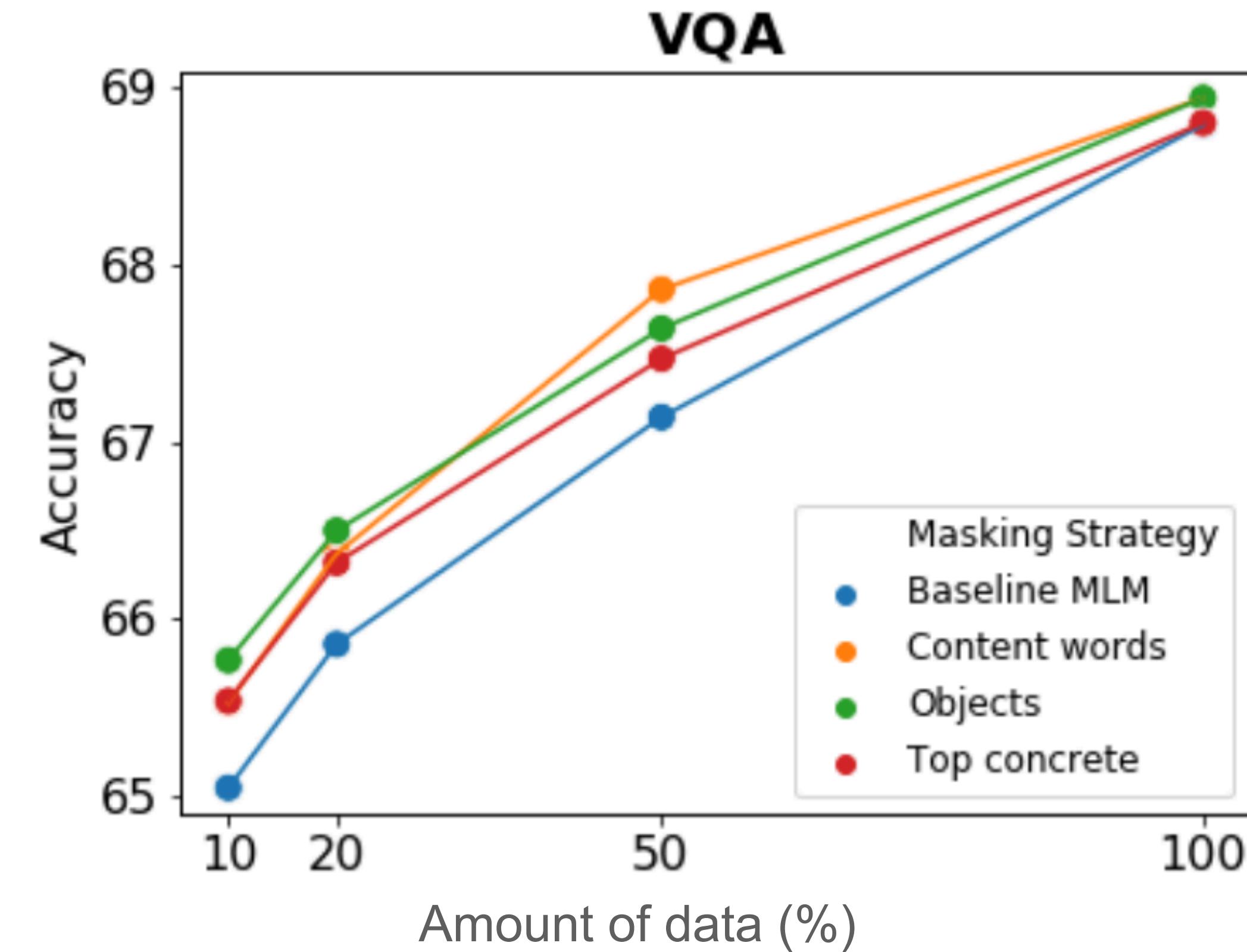


Problems with Pre-training Vision and Language Models

- Of the masked tokens, roughly one half are *stop-words* or *punctuation*
- Image captions are typically short (6-7 words)
 - ⇒ in roughly 1/3 of the captions, **no word is masked!**
- **Improved** pre-training:
 - Mask **exactly one token** in each document
 - Only mask **content words**

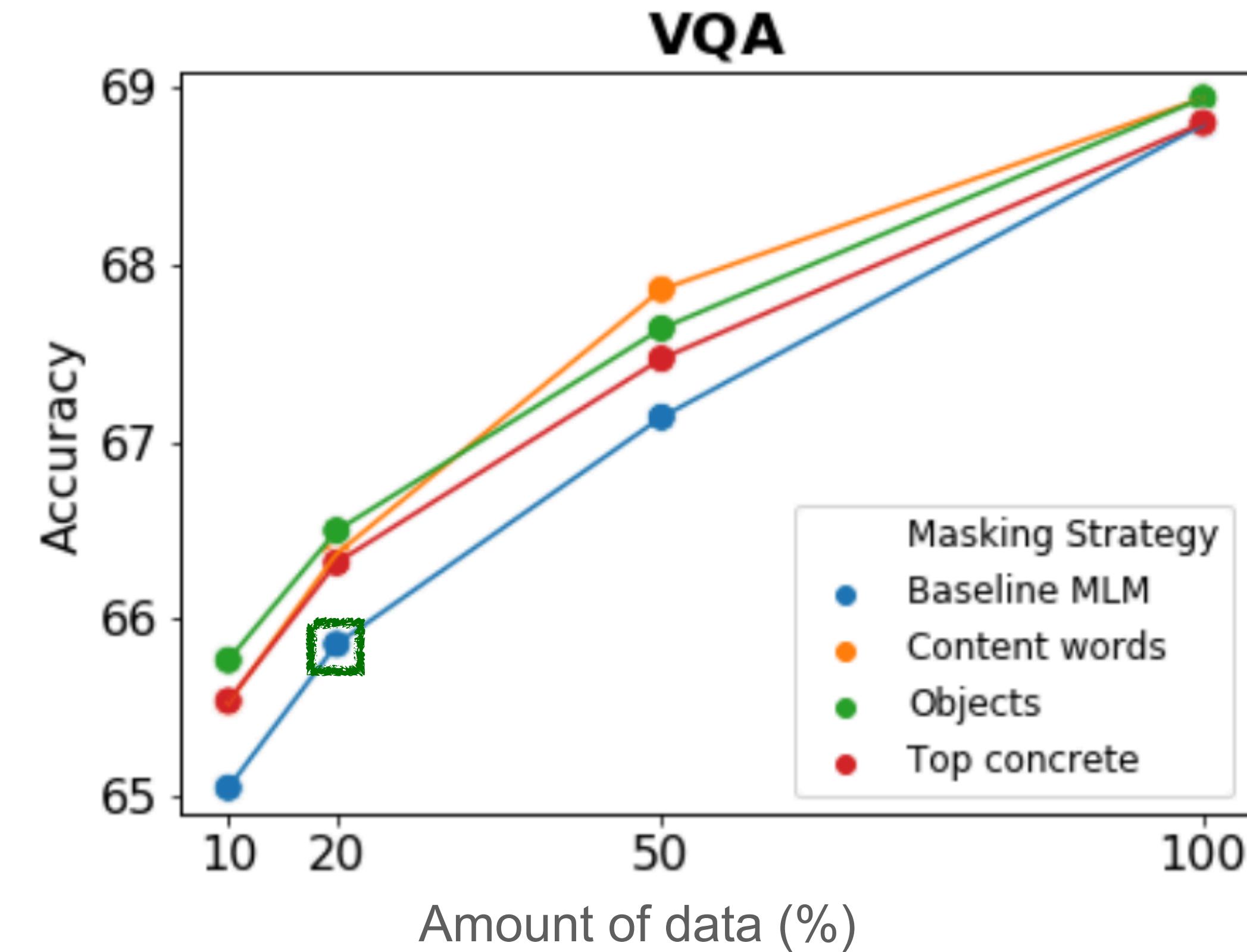


Data Efficient Training



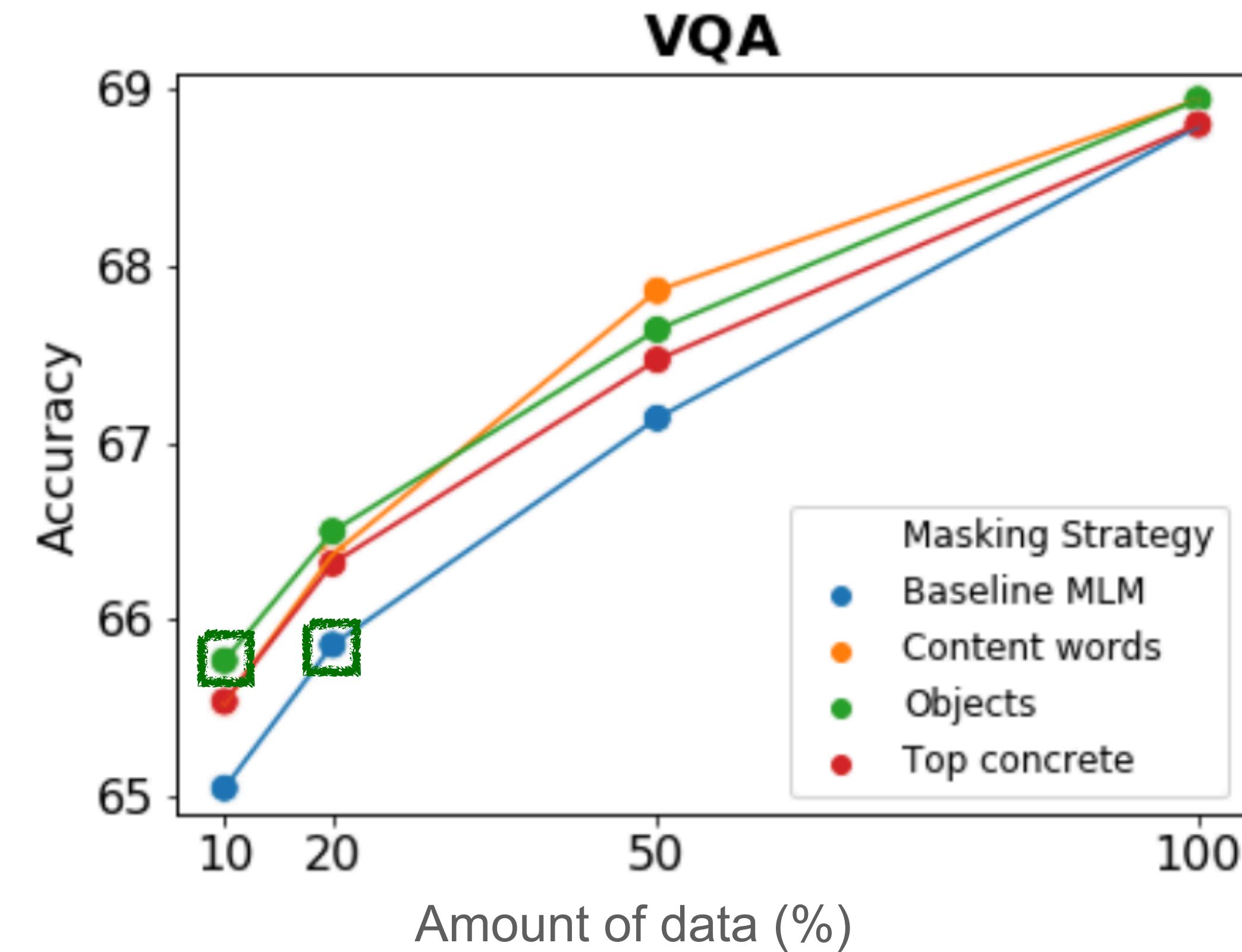


Data Efficient Training



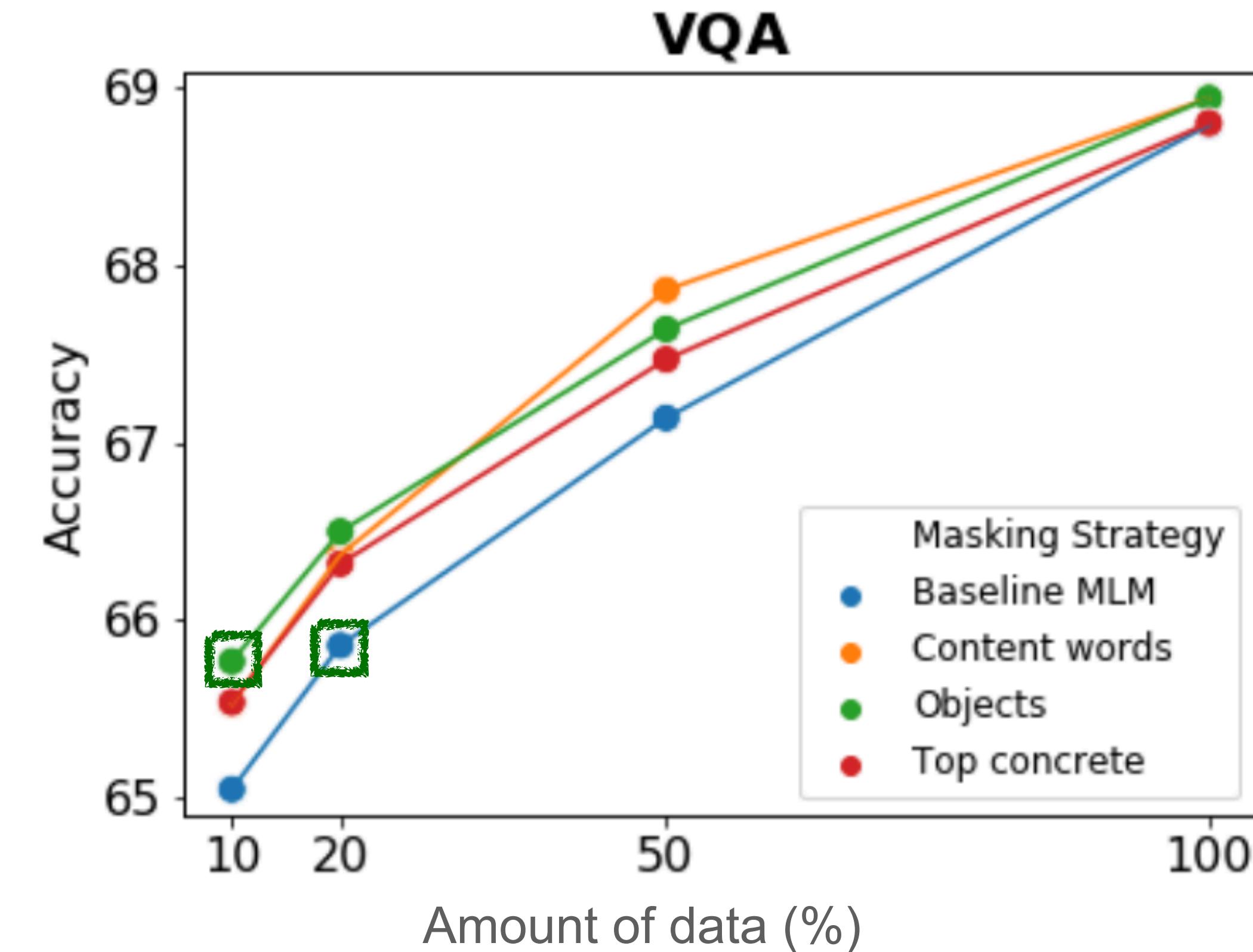


Data Efficient Training





Data Efficient Training



Similar accuracy, twice as fast



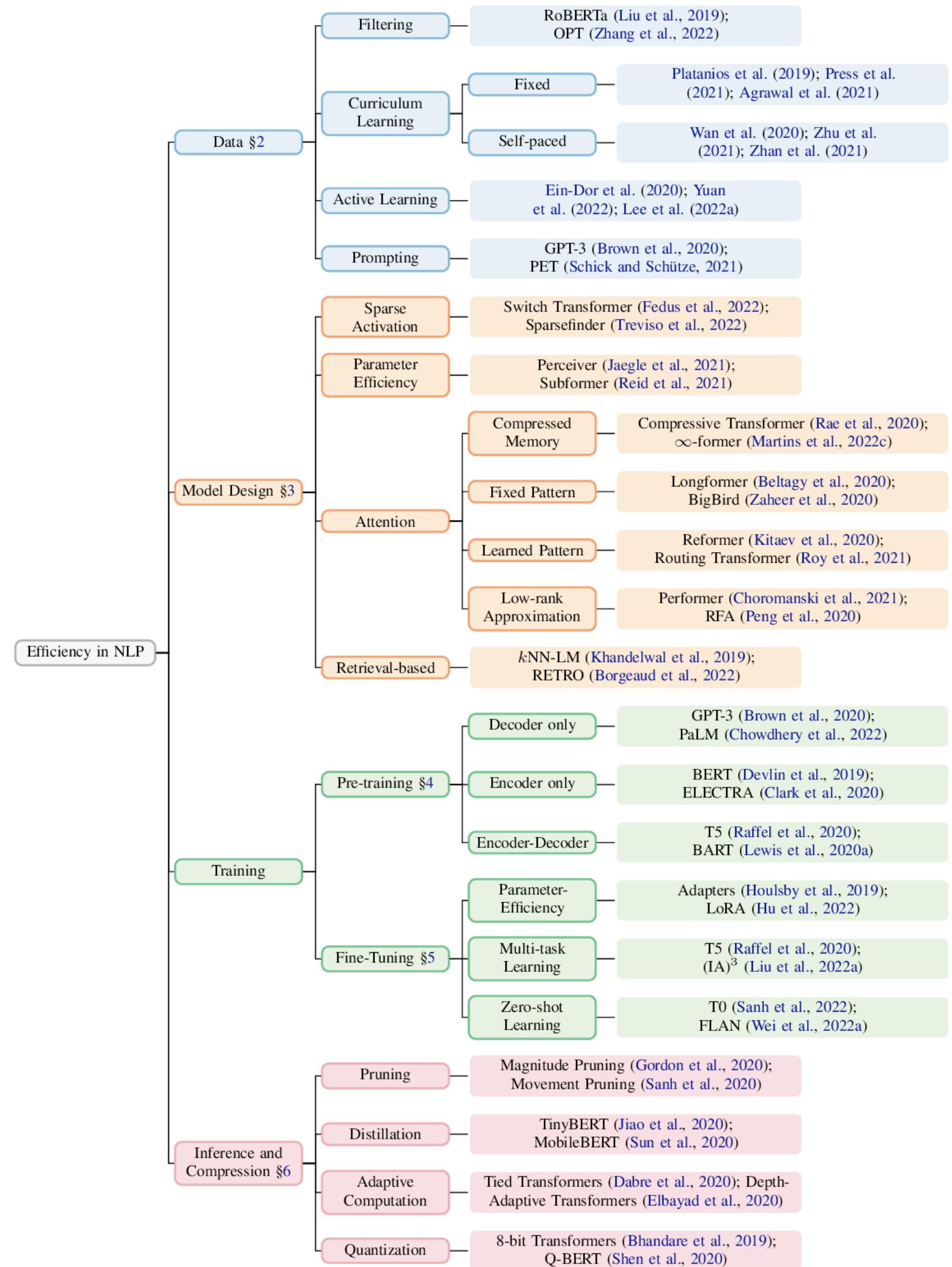
Data Efficiency

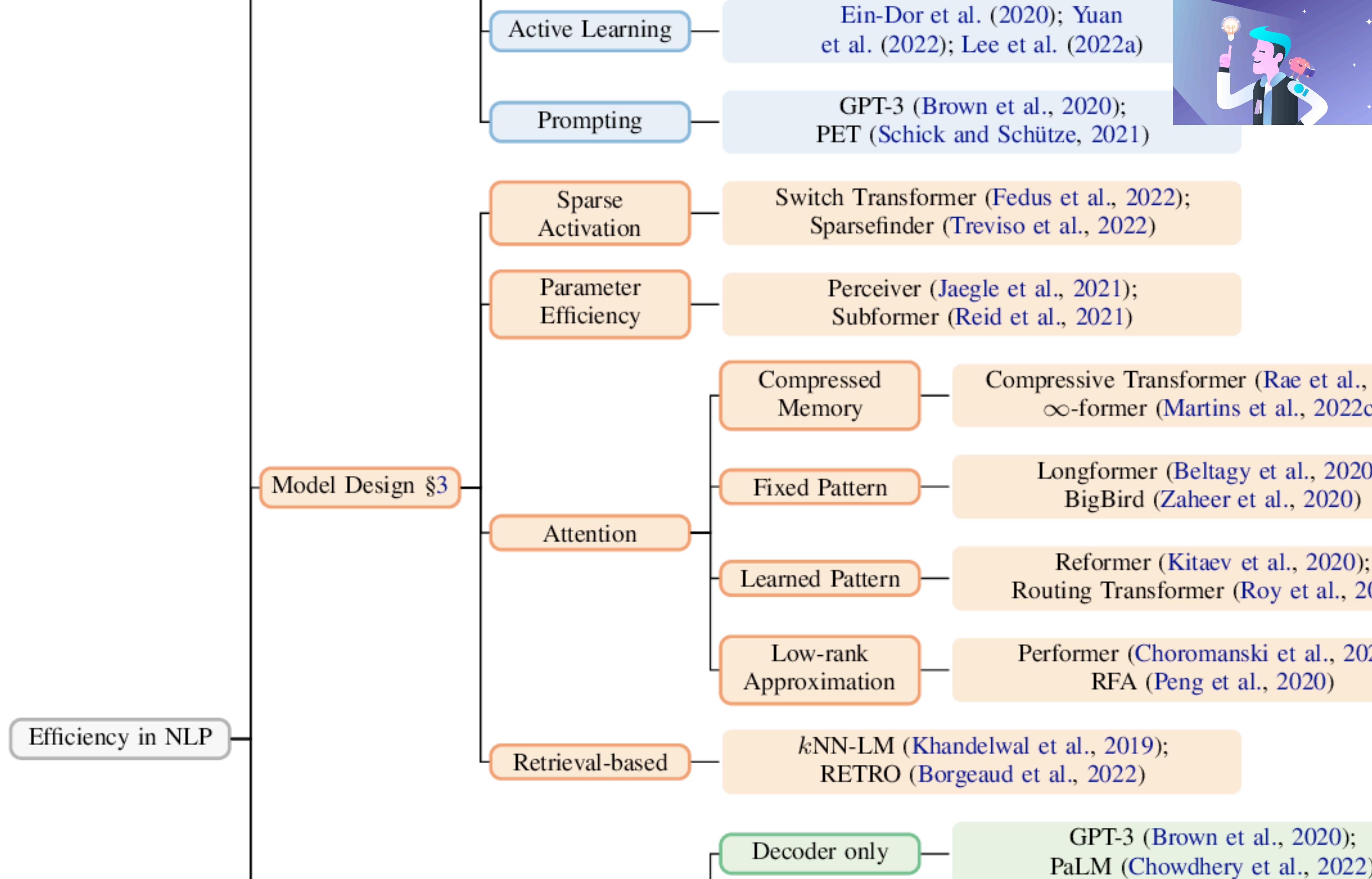
Open Questions

- Do we really need massive web-scale data to train our models?
 - Can we get along with less?



Efficient Methods in NLP





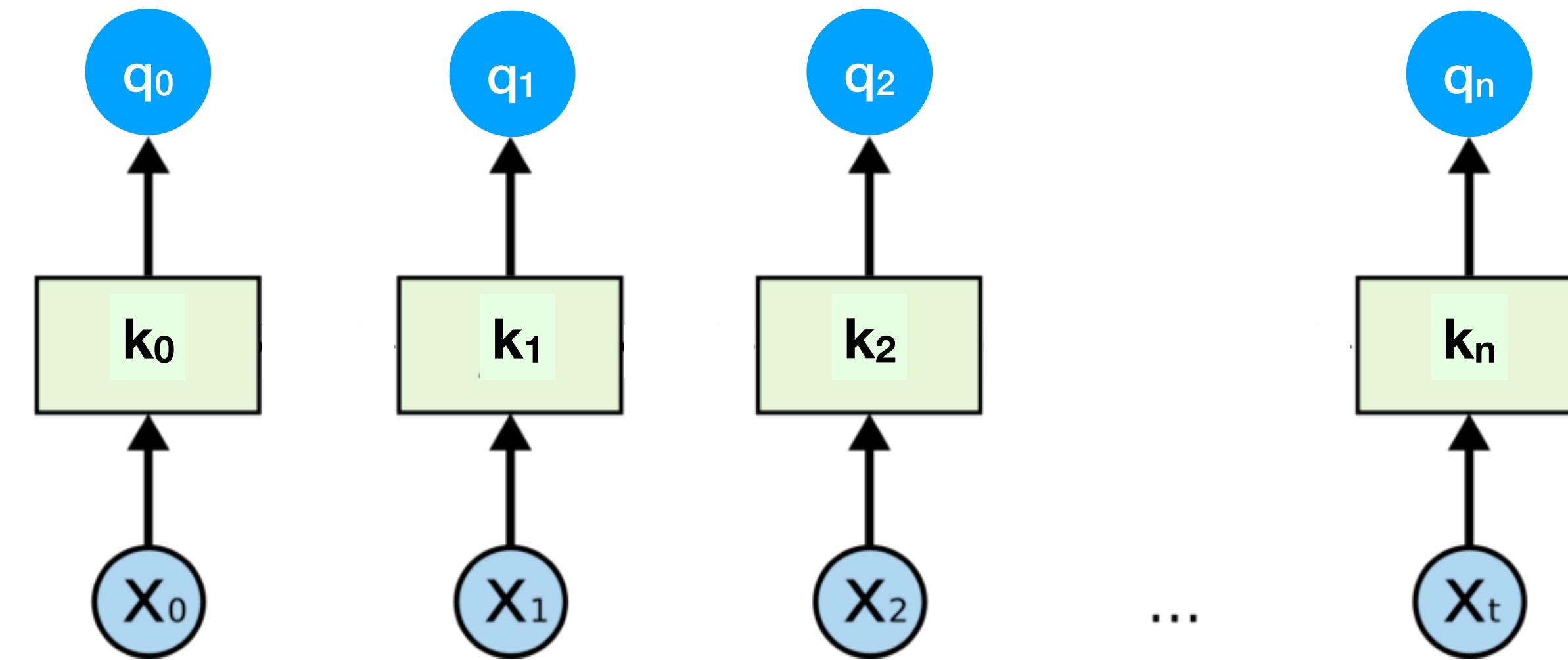


Transformers



Vaswani et al., 2017

- The method for text representation
 - Also for vision, speech, combo, ...



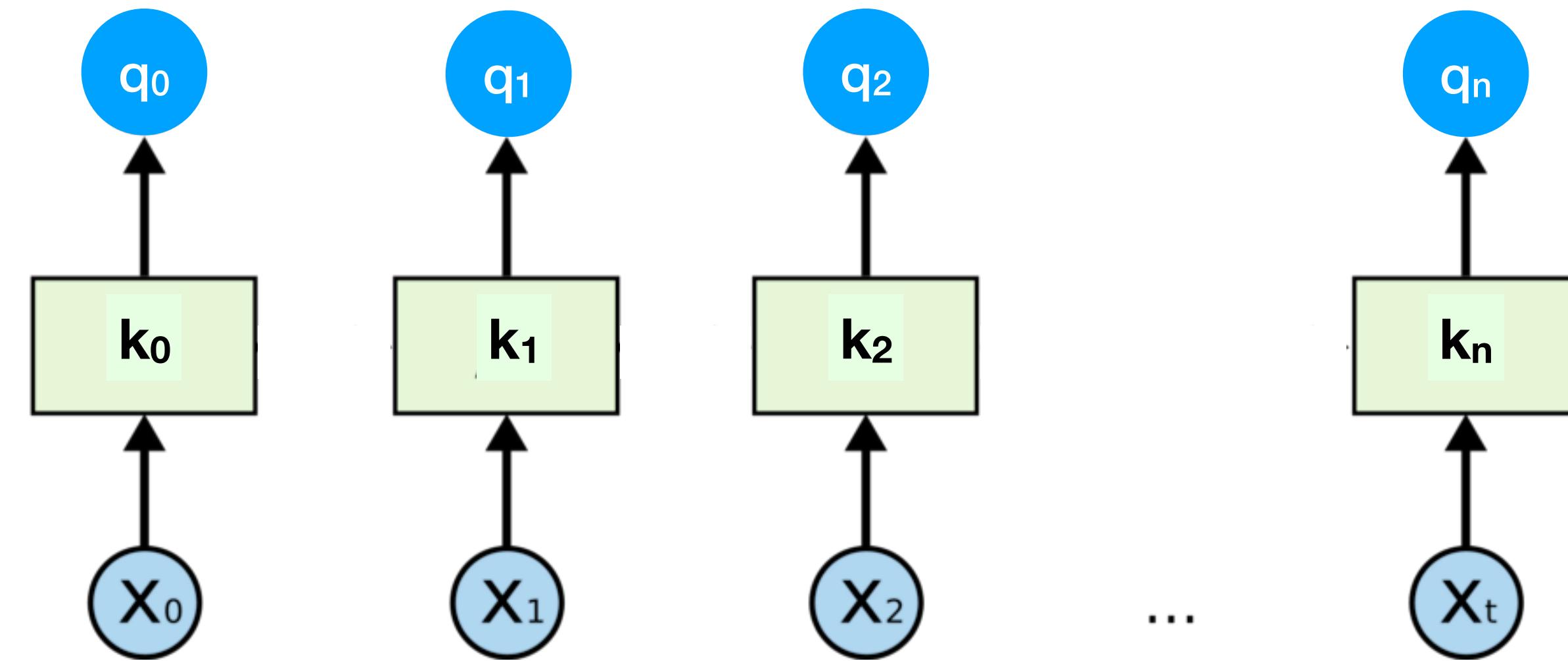


Transformers



Vaswani et al., 2017

- **The** method for text representation
 - Also for vision, speech, combo, ...
- Each word attends to all other words



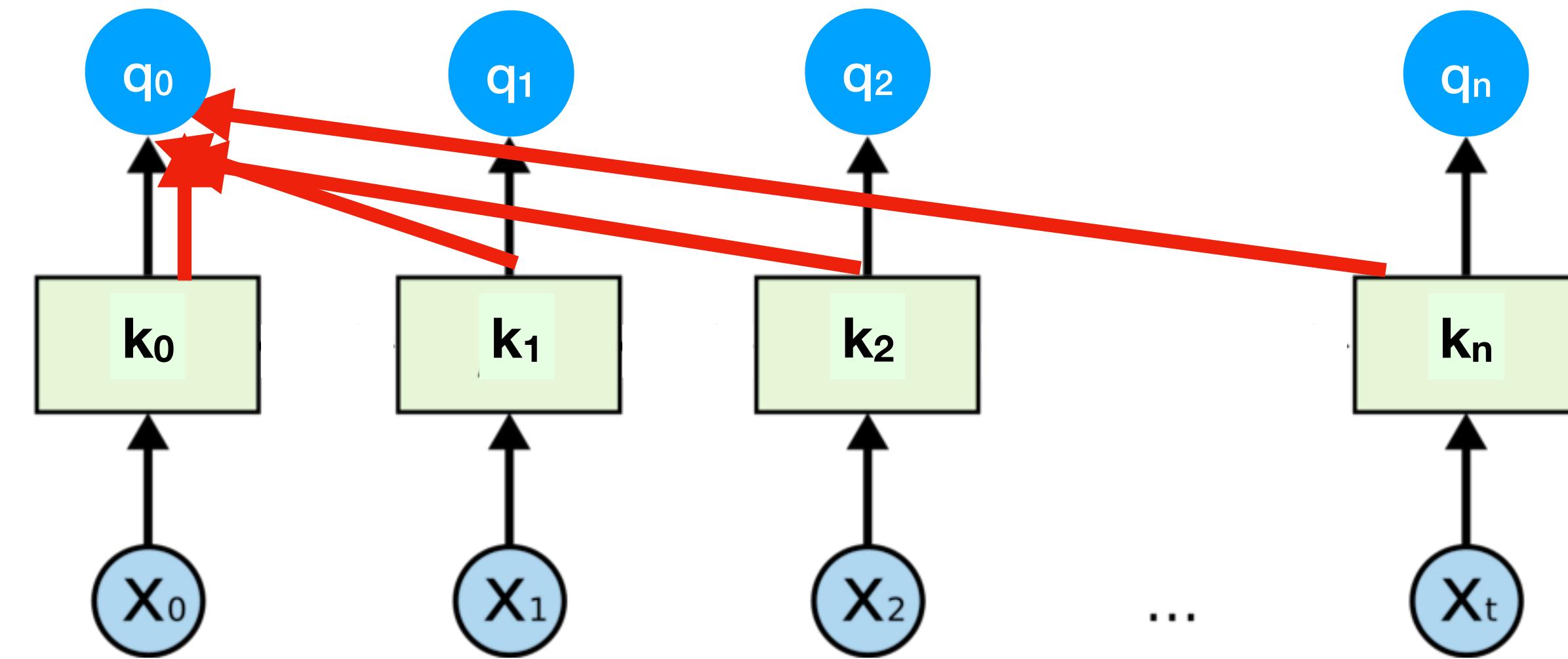


Transformers



Vaswani et al., 2017

- The method for text representation
 - Also for vision, speech, combo, ...
- Each word attends to all other words



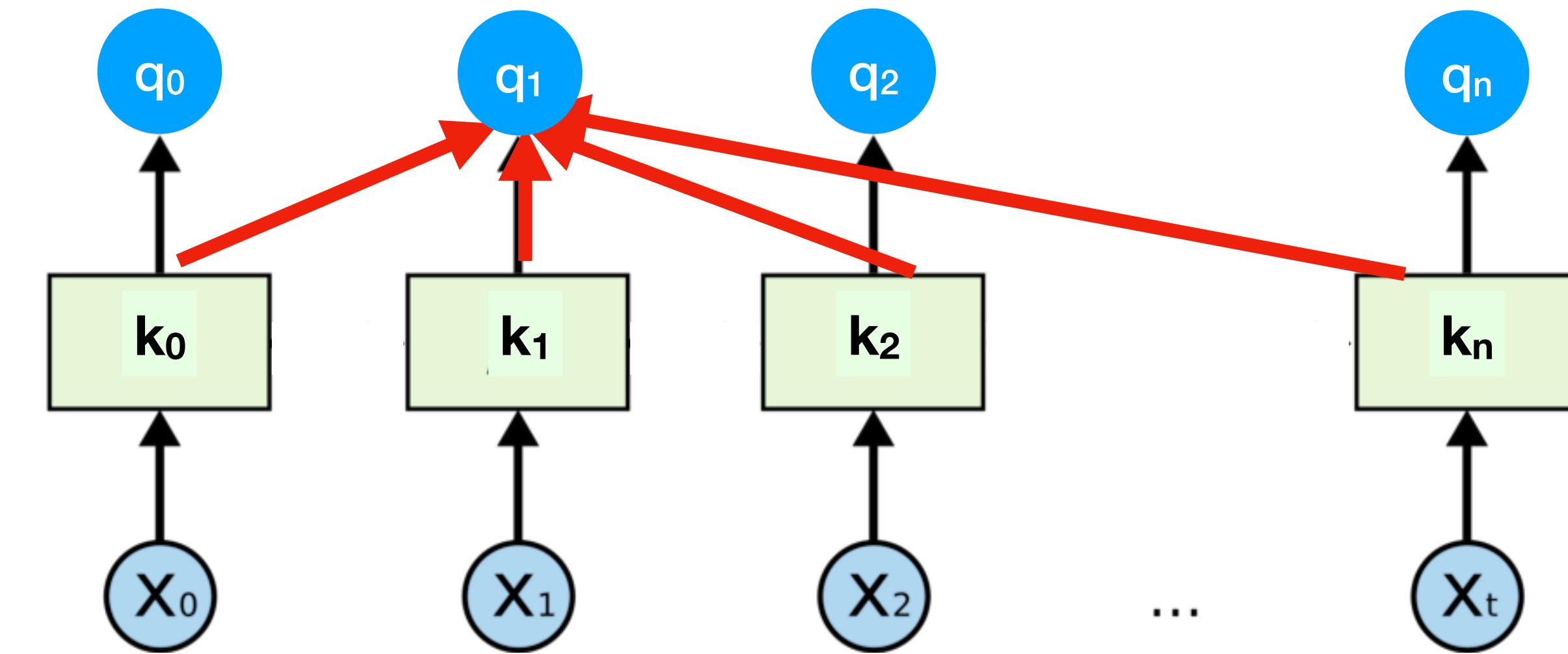


Transformers



Vaswani et al., 2017

- **The** method for text representation
 - Also for vision, speech, combo, ...
- Each word attends to all other words



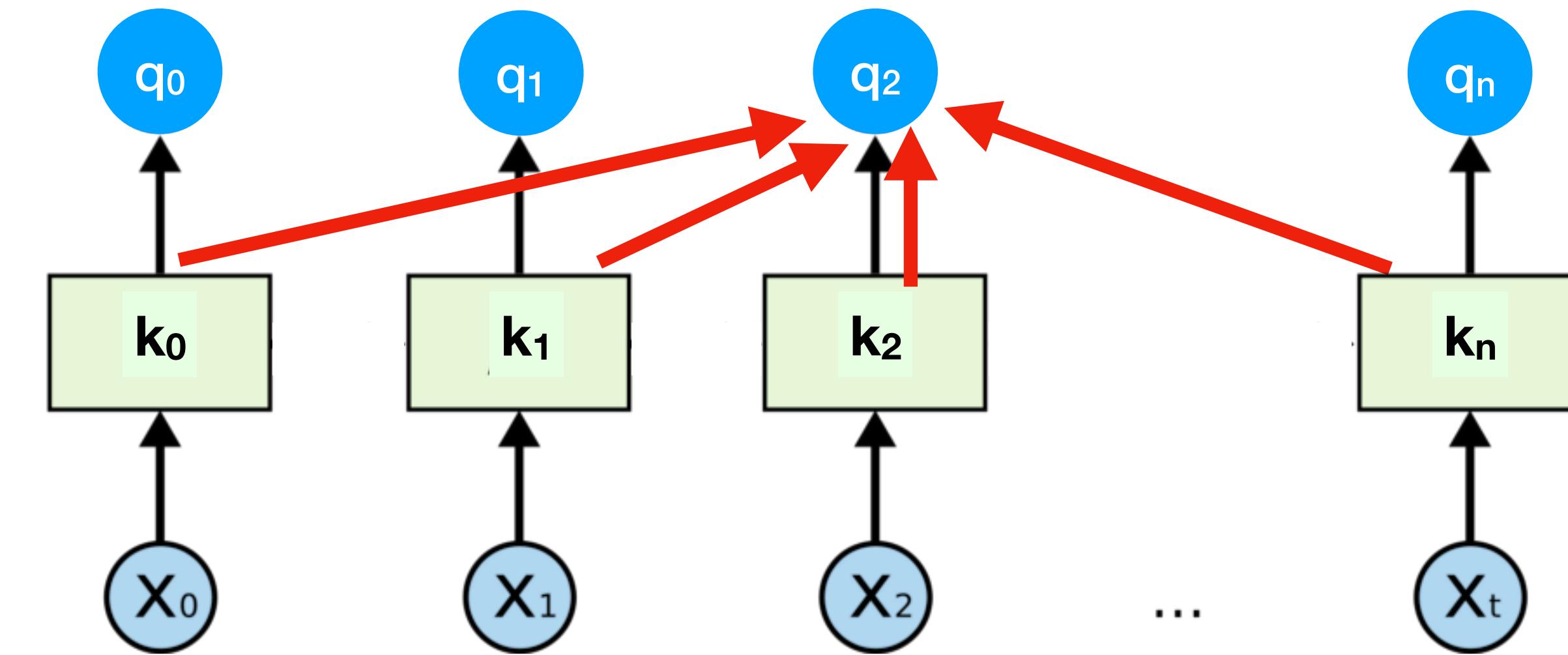


Transformers



Vaswani et al., 2017

- **The** method for text representation
 - Also for vision, speech, combo, ...
- Each word attends to all other words



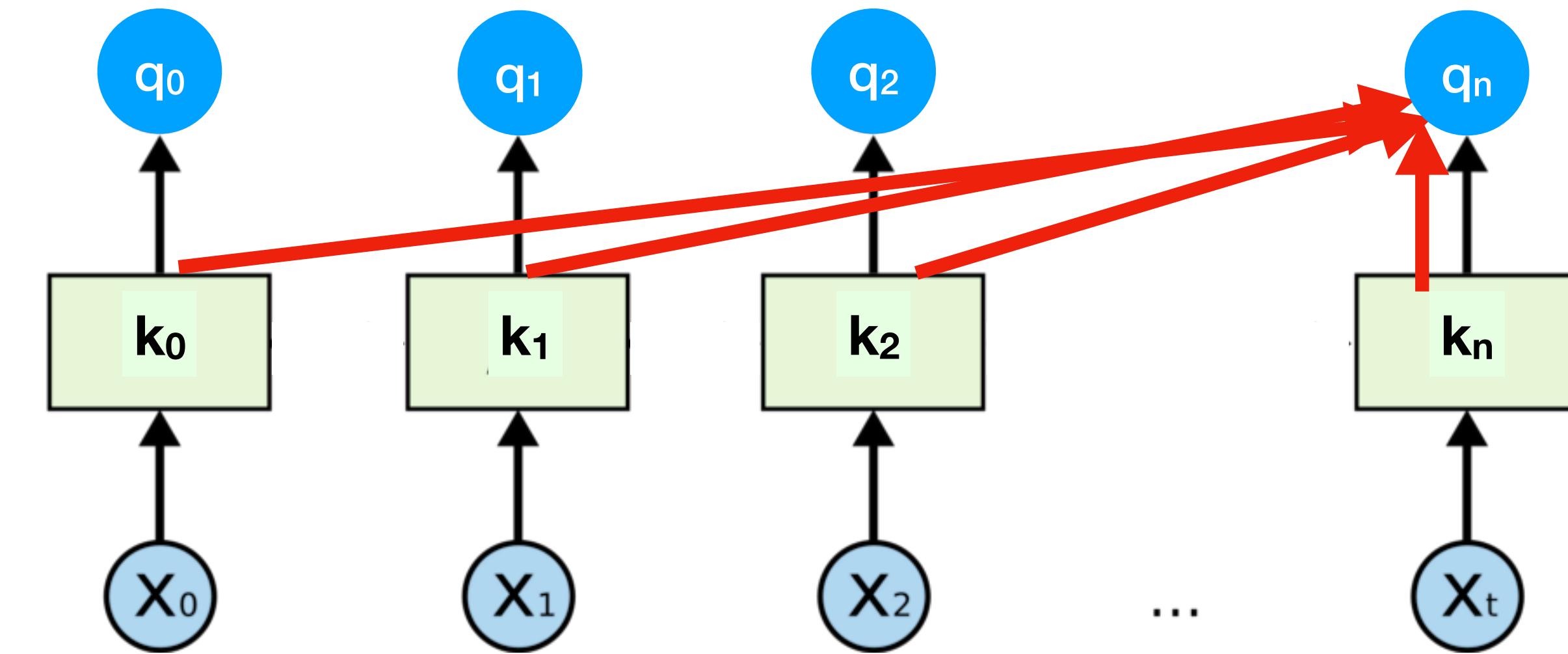


Transformers



Vaswani et al., 2017

- **The** method for text representation
 - Also for vision, speech, combo, ...
- Each word attends to all other words



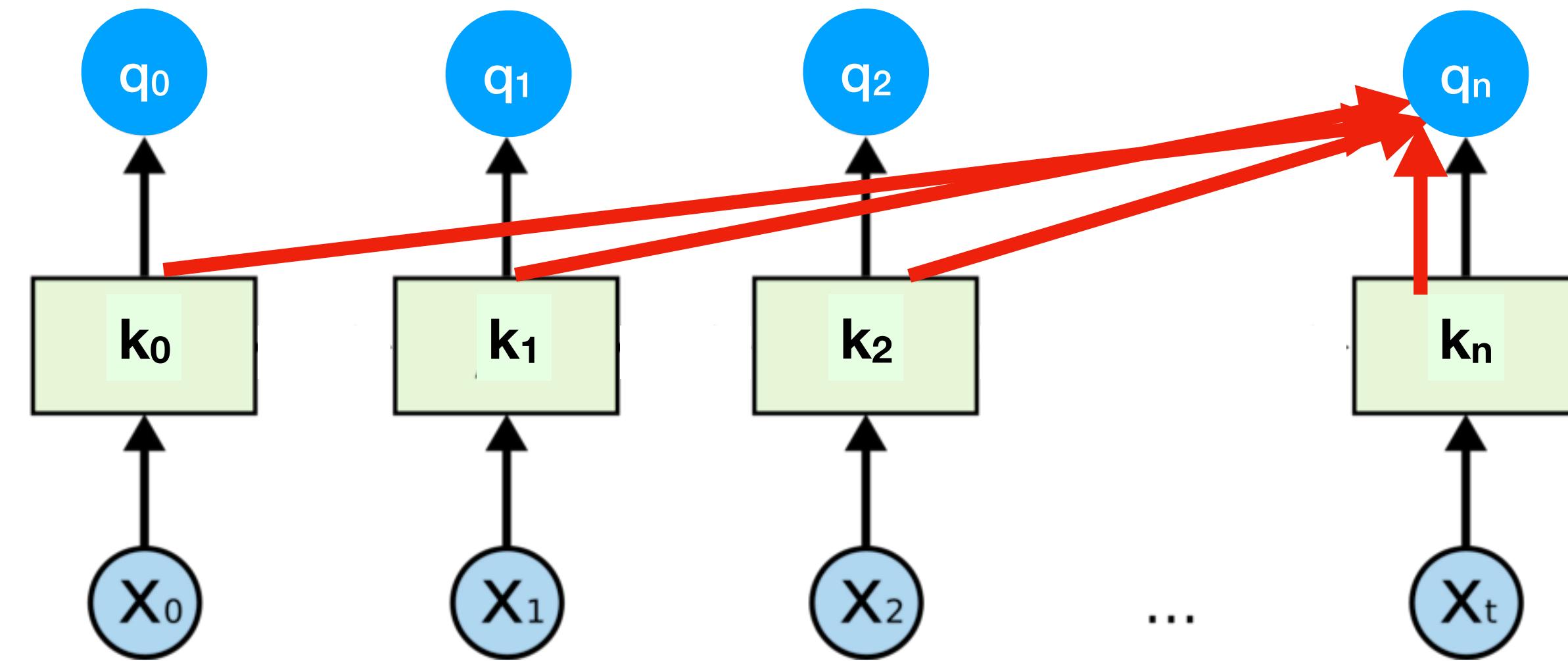


Transformers



Vaswani et al., 2017

- **The** method for text representation
 - Also for vision, speech, combo, ...
- Each word attends to all other words
- $O(n^2)$ complexity in the sentence length n



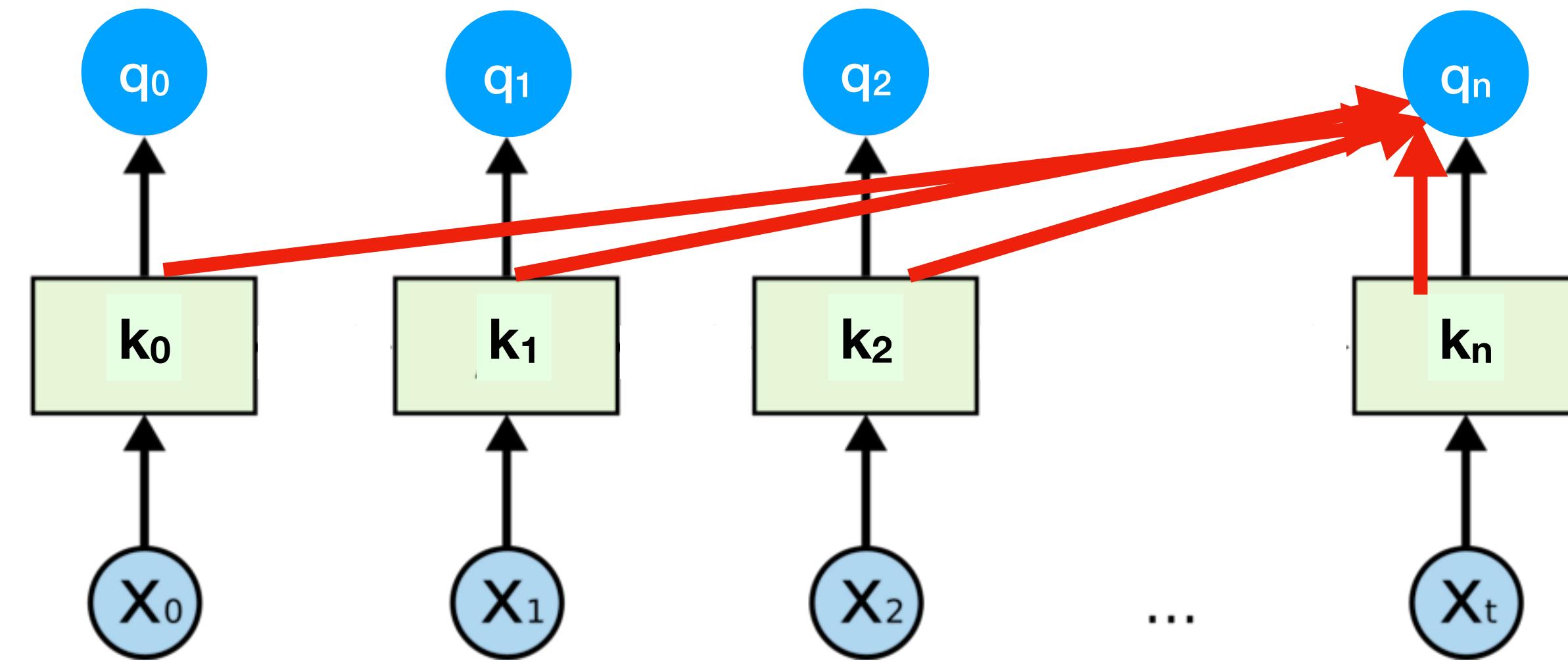


Transformers



Vaswani et al., 2017

- **The** method for text representation
 - Also for vision, speech, combo, ...
- Each word attends to all other words
- $O(n^2)$ complexity in the sentence length n
- Fatal for long sequences
 - Books, articles, etc.



Random Feature Attention

Peng, Pappas, Yogatama, **S.**, Smith, & Kong, ICLR 2021

spotlight presentation

- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)



Random Feature Attention

Peng, Pappas, Yogatama, **S.**, Smith, & Kong, ICLR 2021

spotlight presentation

- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)
 - Some math



Random Feature Attention

Peng, Pappas, Yogatama, **S.**, Smith, & Kong, ICLR 2021

spotlight presentation

- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)
 - Some math



Random Feature Attention

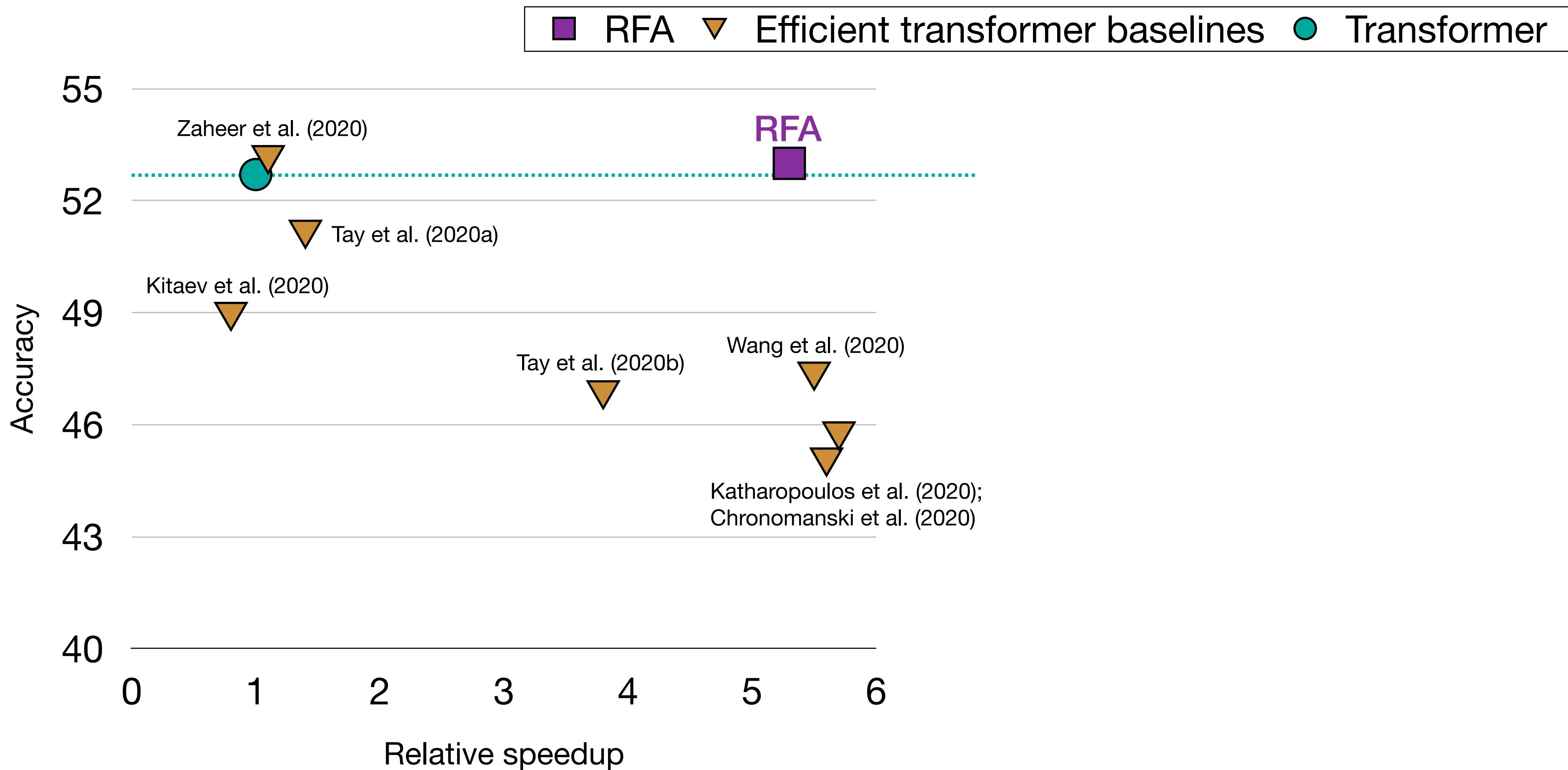
Peng, Pappas, Yogatama, **S.**, Smith, & Kong, ICLR 2021

spotlight presentation

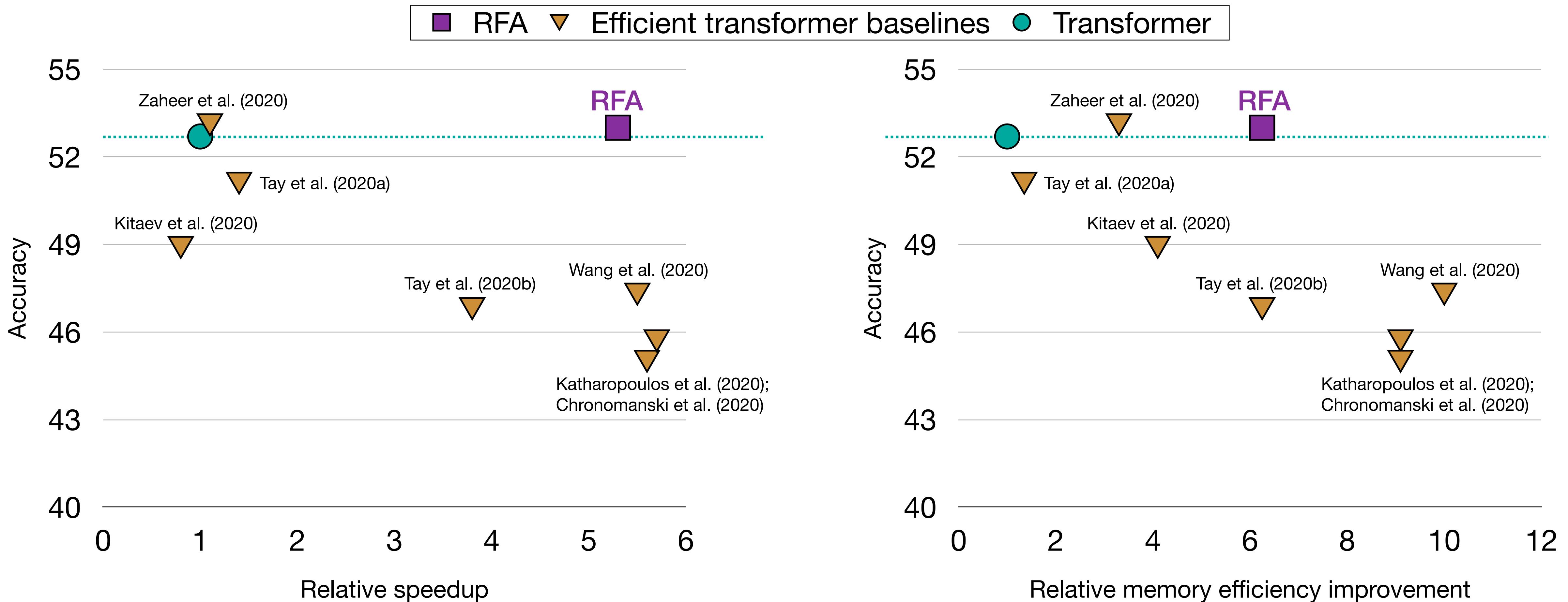
- **Key idea:** approximate the attention function using random Fourier features
 - Rahimi and Recht (2007)
- Some math
- Linear runtime and memory requirements



Better Efficiency-Accuracy Tradeoff



Better Efficiency-Accuracy Tradeoff

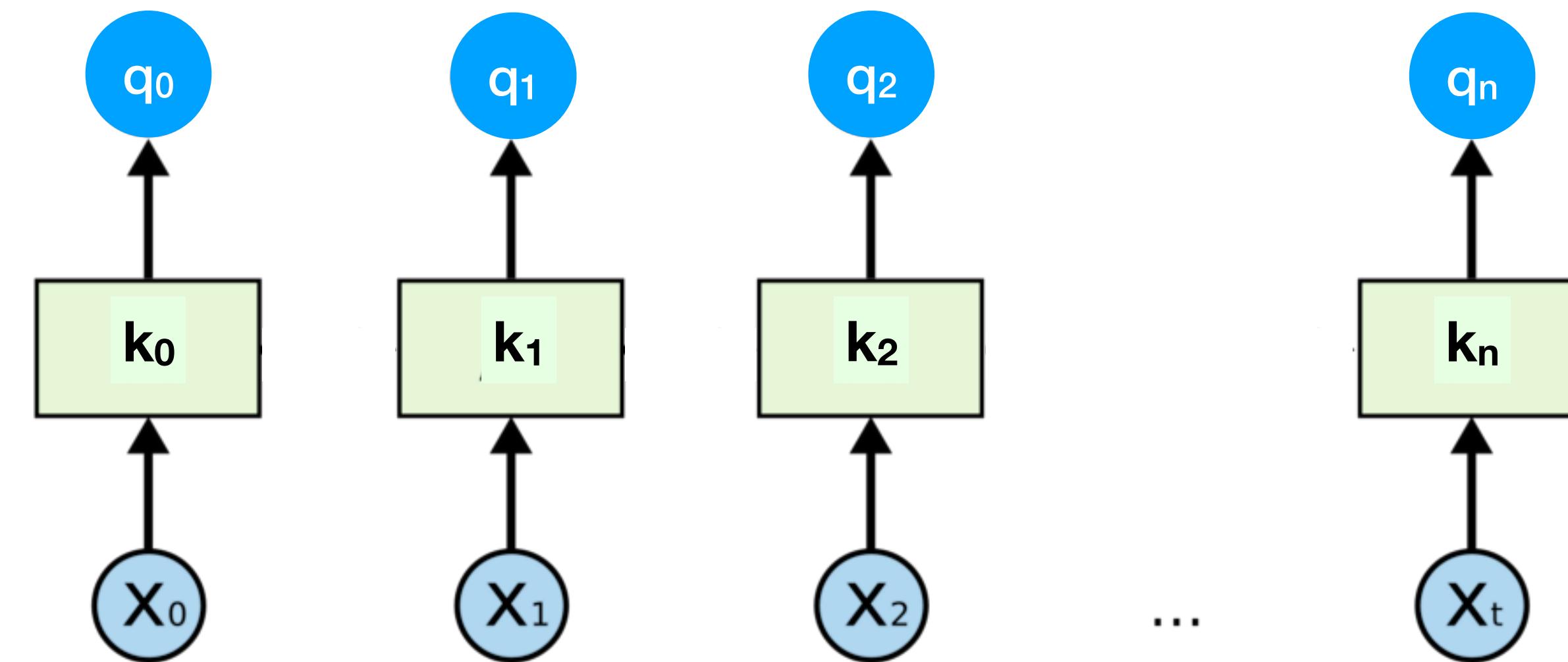




ABC: Attention with Bounded-memory Control

Peng, Kasai, Pappas, Yogatama, Wu, Kong, S. & Smith, ACL 2022

- **Key intuition:** treat the sentence as **memory of size n**

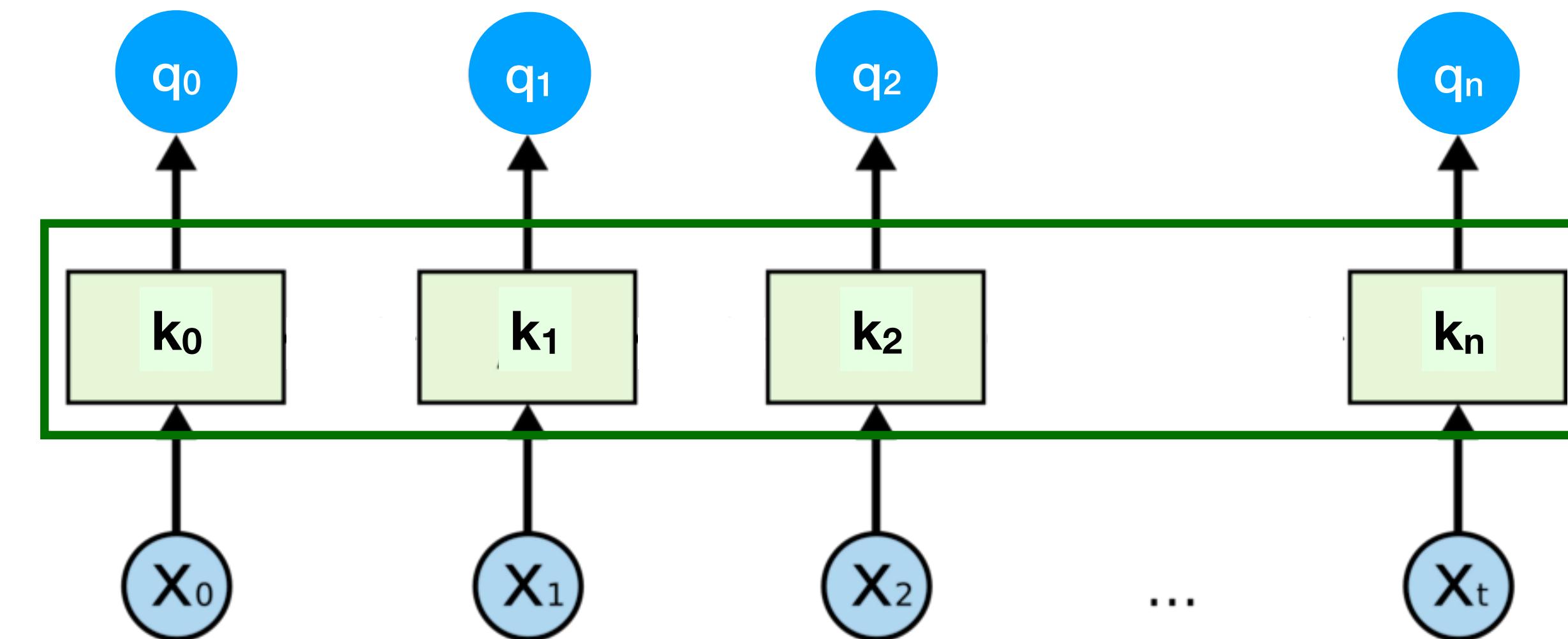




ABC: Attention with Bounded-memory Control

Peng, Kasai, Pappas, Yogatama, Wu, Kong, S. & Smith, ACL 2022

- **Key intuition:** treat the sentence as **memory of size n**

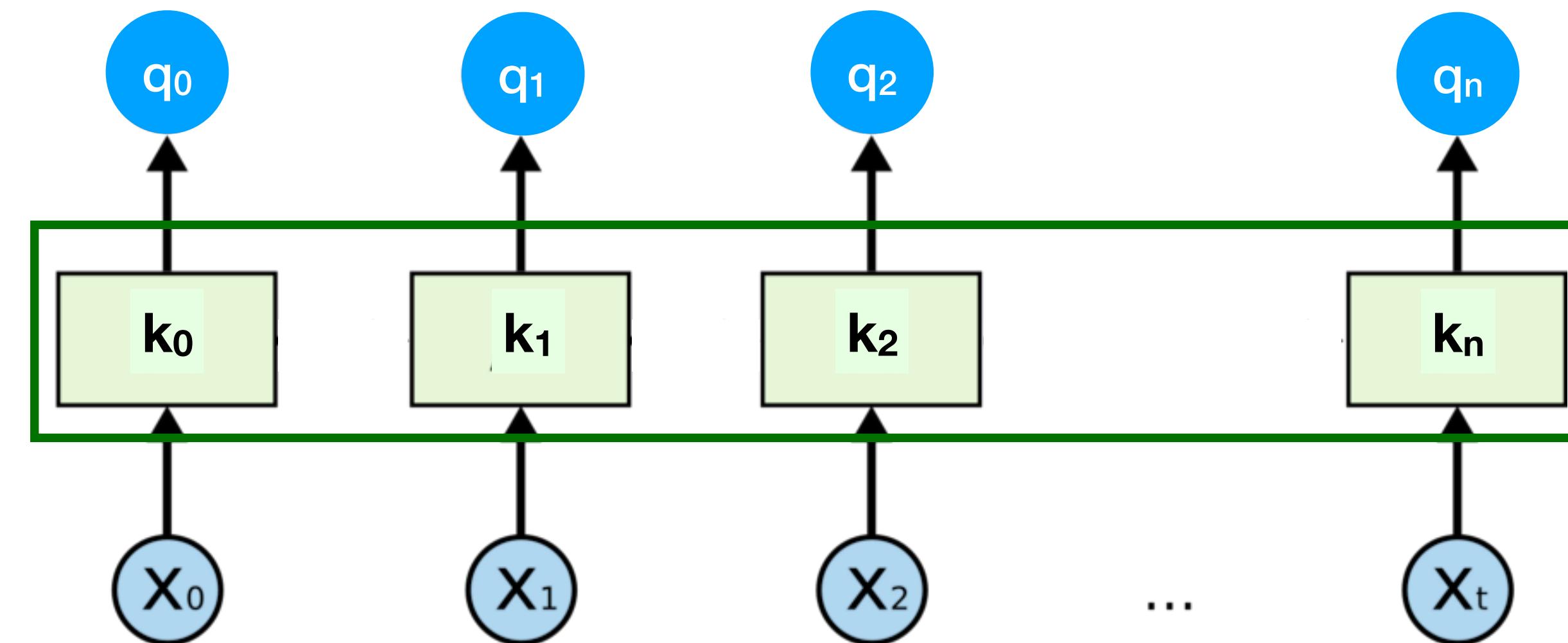




ABC: Attention with Bounded-memory Control

Peng, Kasai, Pappas, Yogatama, Wu, Kong, **S.** & Smith, ACL 2022

- **Key intuition:** treat the sentence as **memory of size n**
- **Key idea:** replace this memory with a fixed size memory of (fixed) size $k \ll n$
 - Instead of attending n tokens, each word attends to k tokens

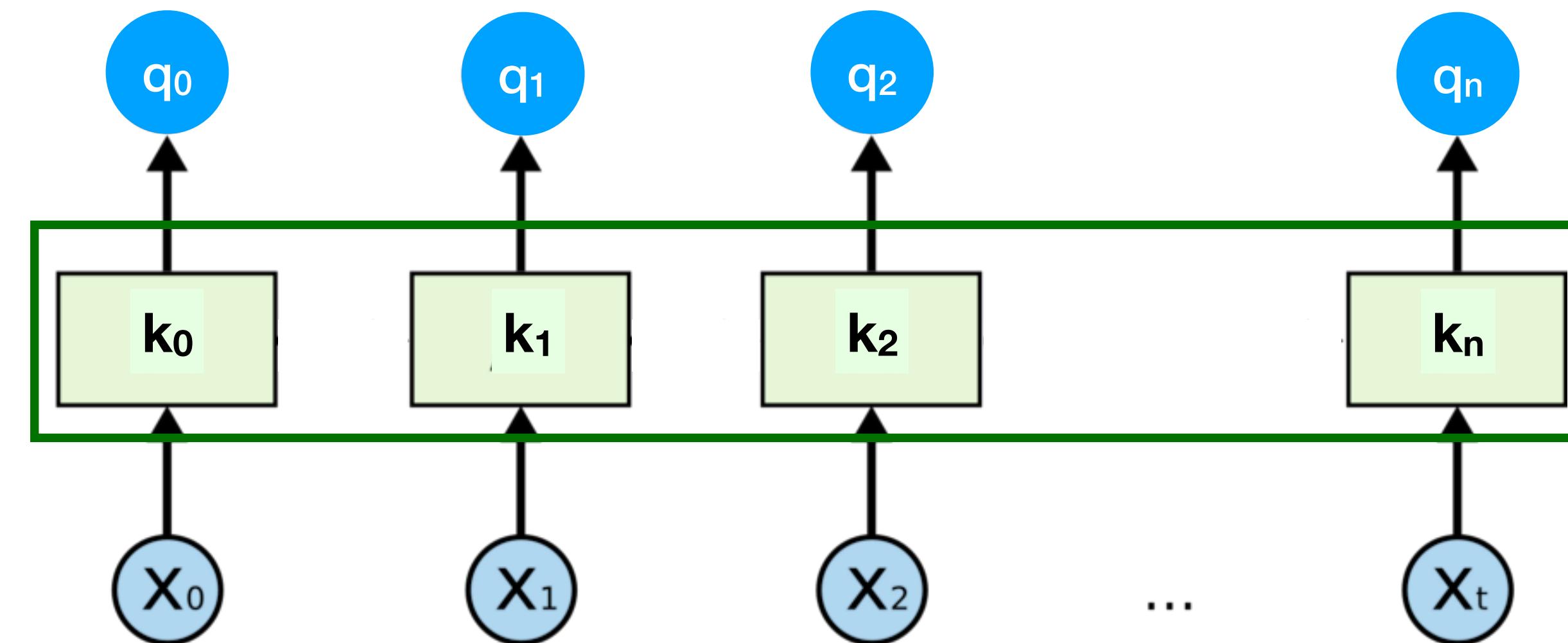




ABC: Attention with Bounded-memory Control

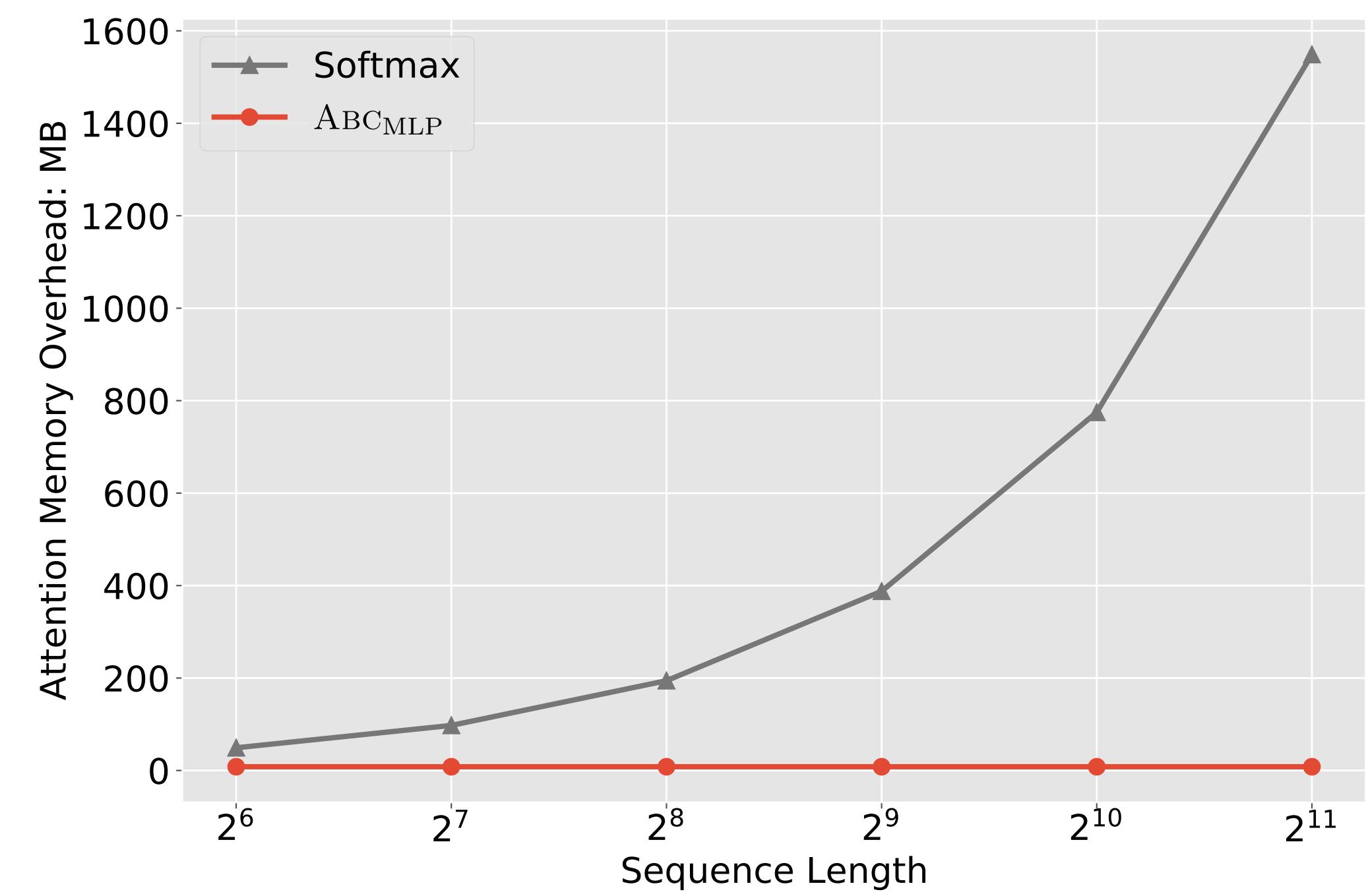
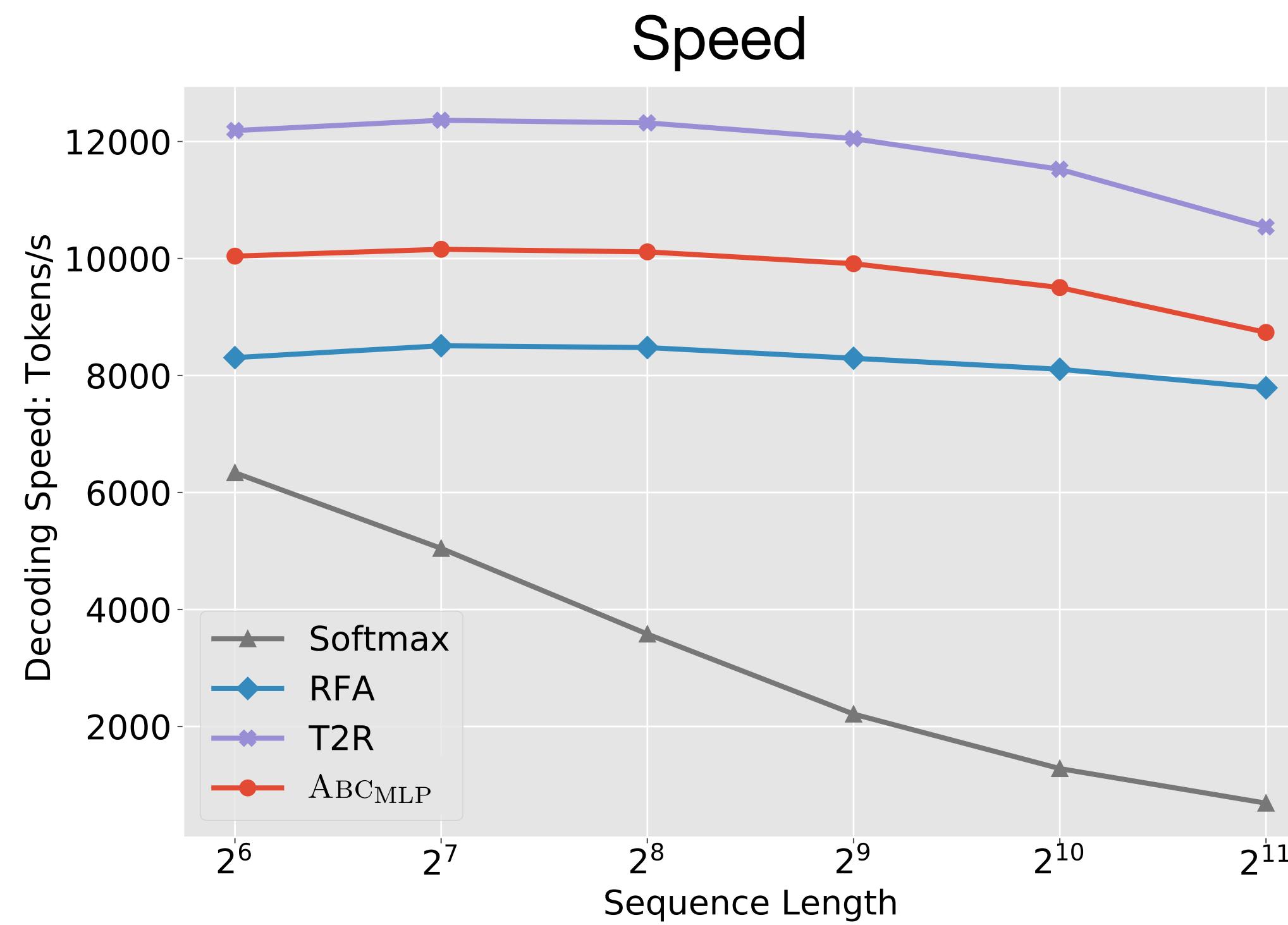
Peng, Kasai, Pappas, Yogatama, Wu, Kong, S. & Smith, ACL 2022

- **Key intuition:** treat the sentence as **memory of size n**
- **Key idea:** replace this memory with a fixed size memory of (fixed) size $k \ll n$
 - Instead of attending n tokens, each word attends to k tokens
- Overall complexity linear in n
 - With constant k





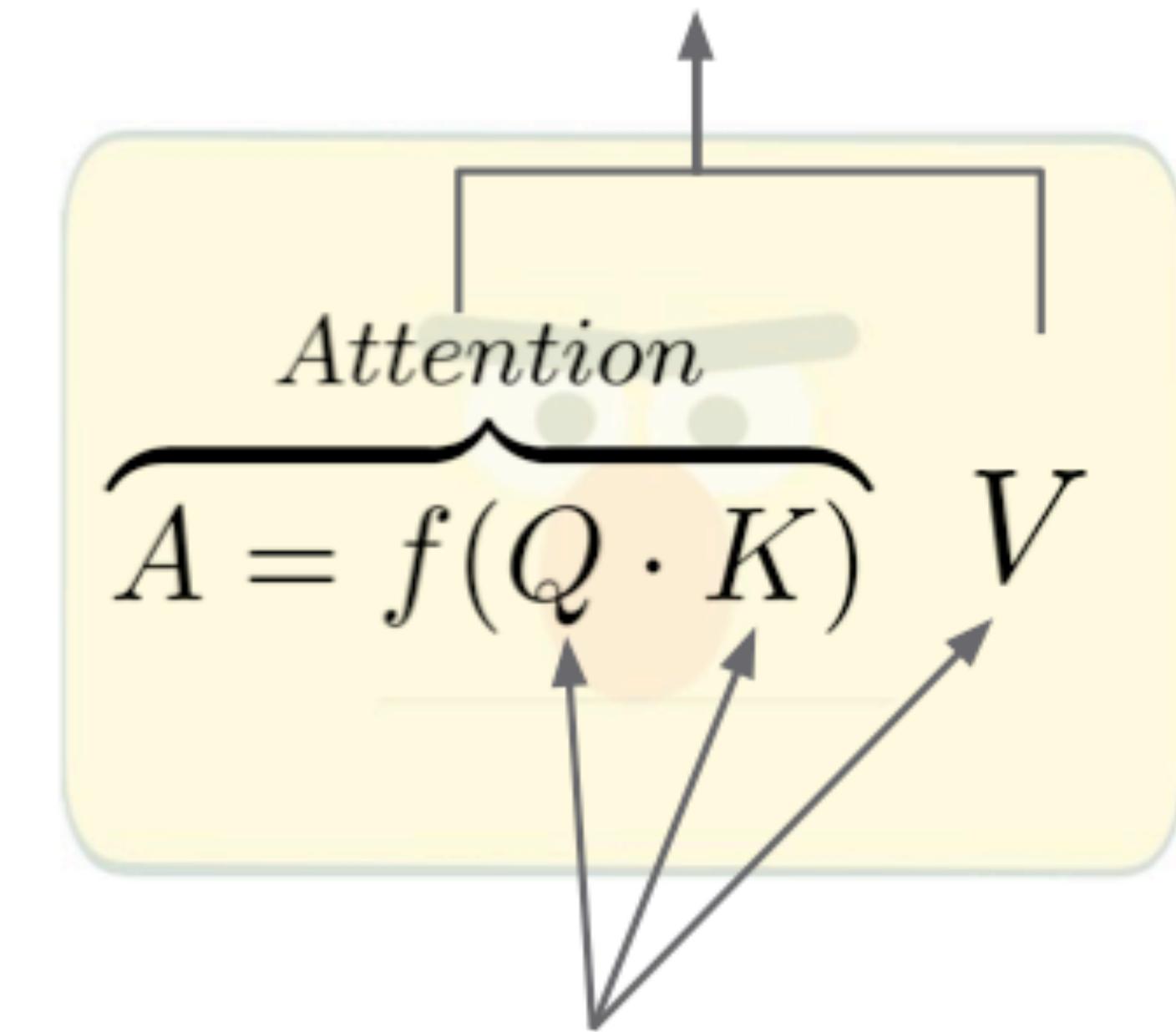
ABC Results



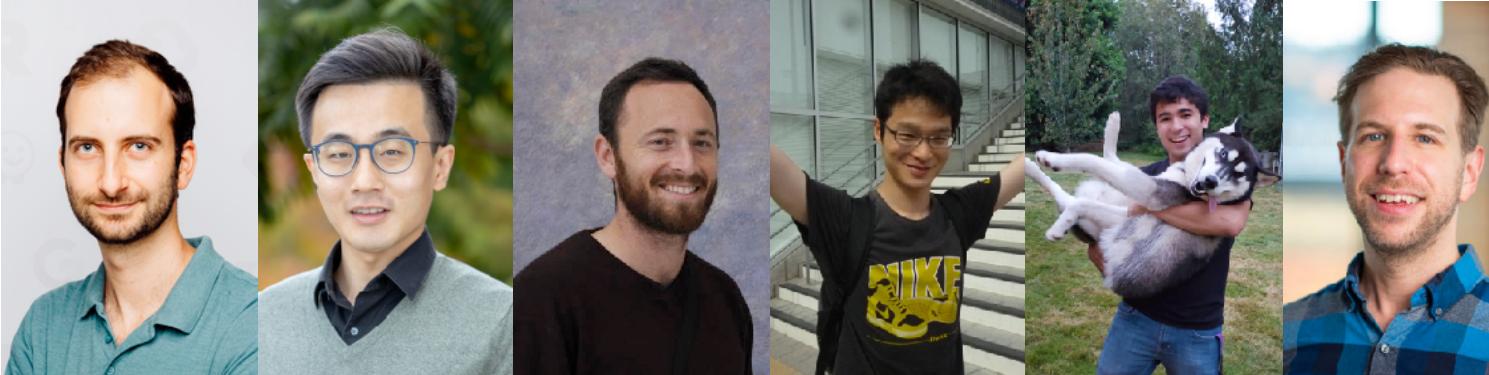


How Much Does Attention Actually Attend?

Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022



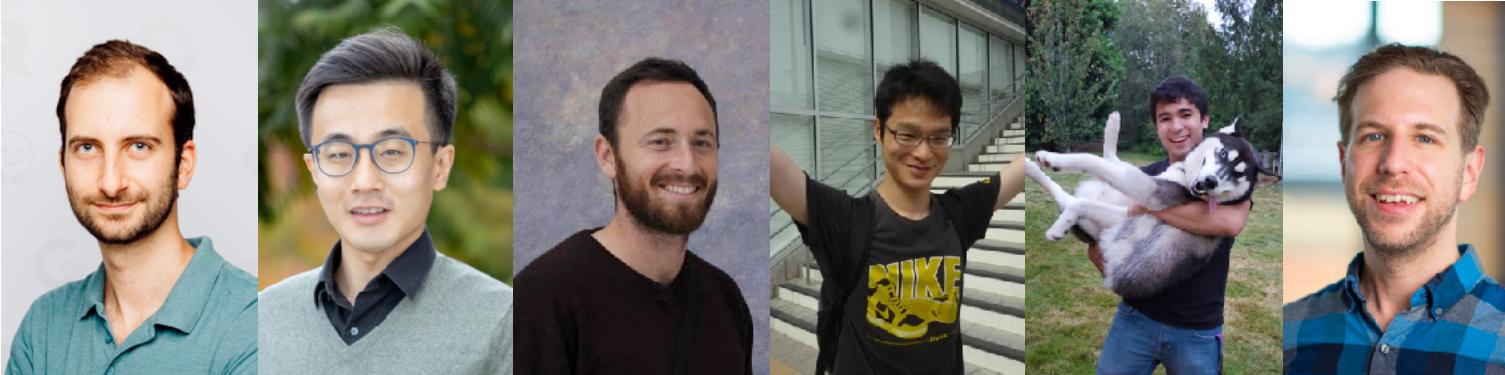
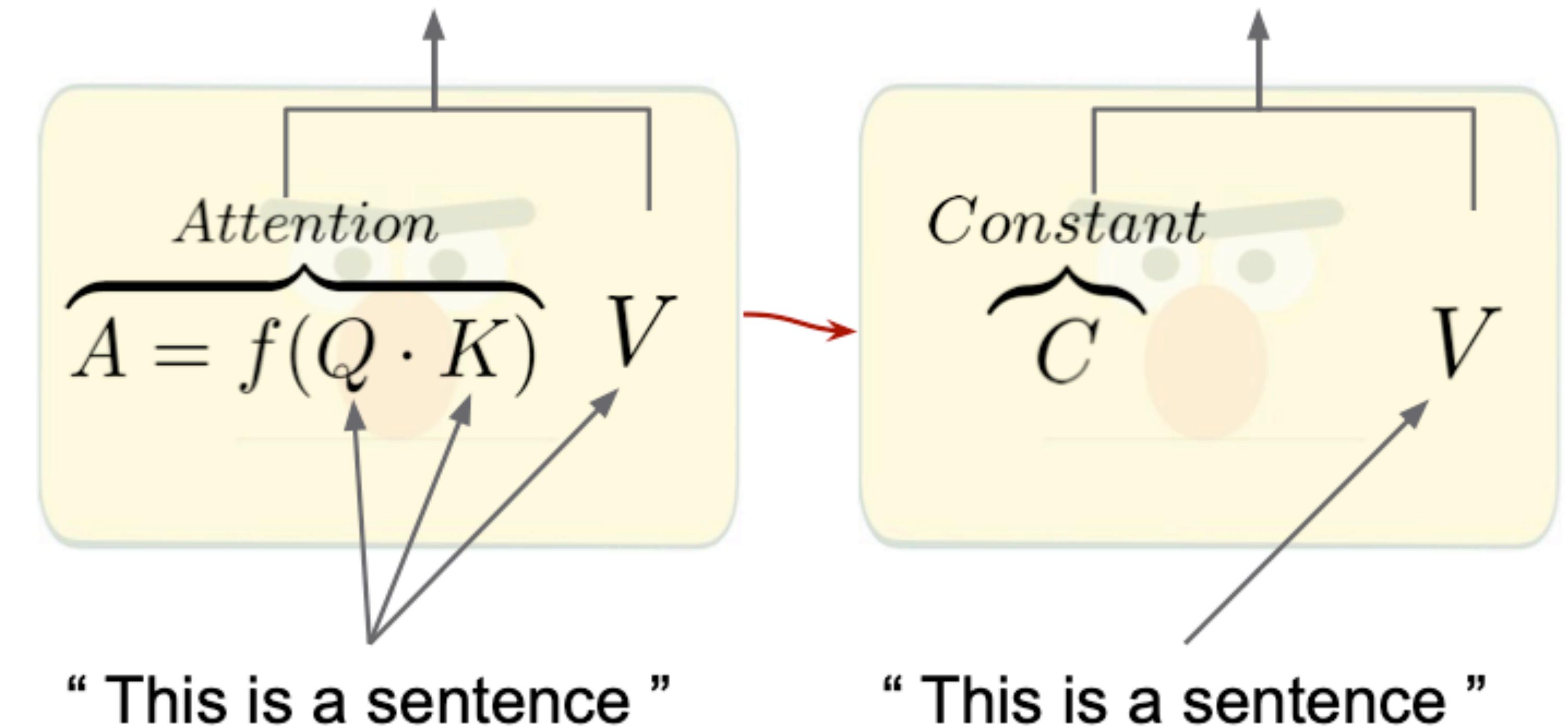
“ This is a sentence ”





How Much Does Attention Actually Attend?

Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022

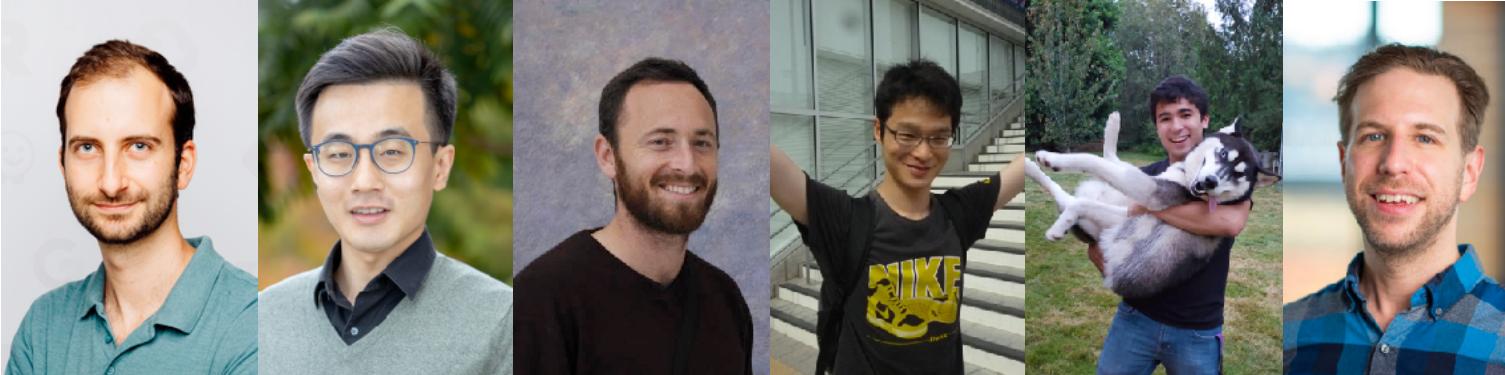
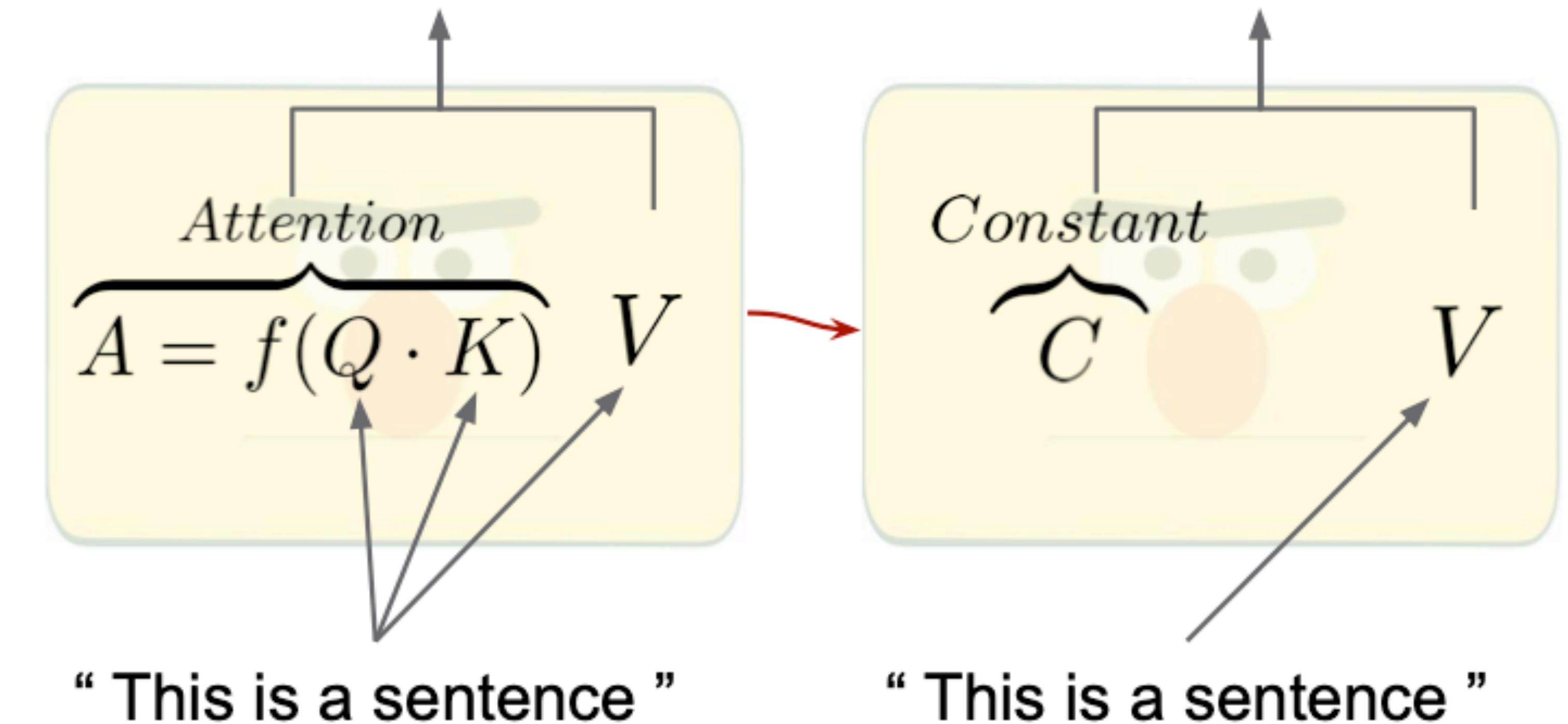




How Much Does Attention Actually Attend?

Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022

- Model doesn't collapse
 - Average accuracy loss of 8% only

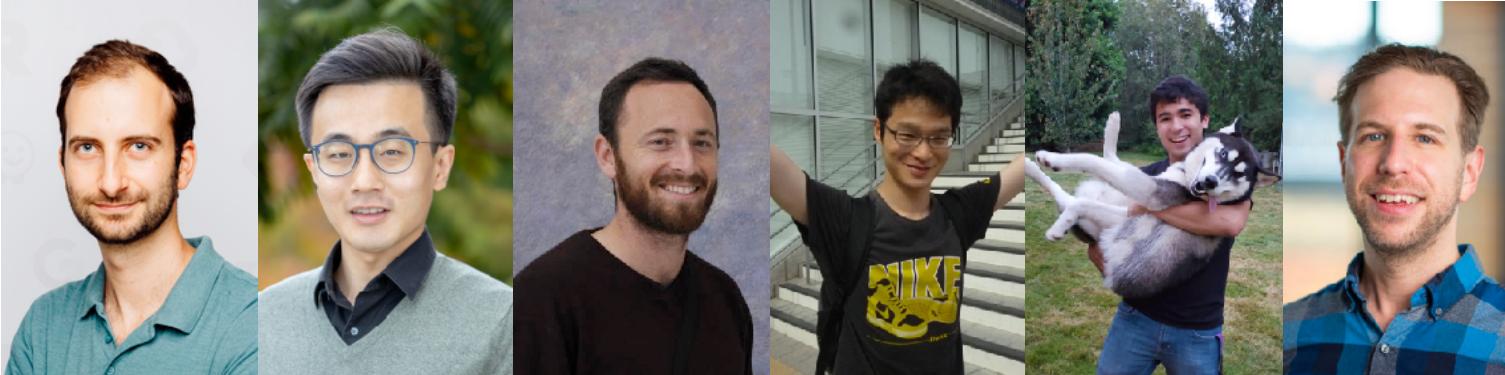
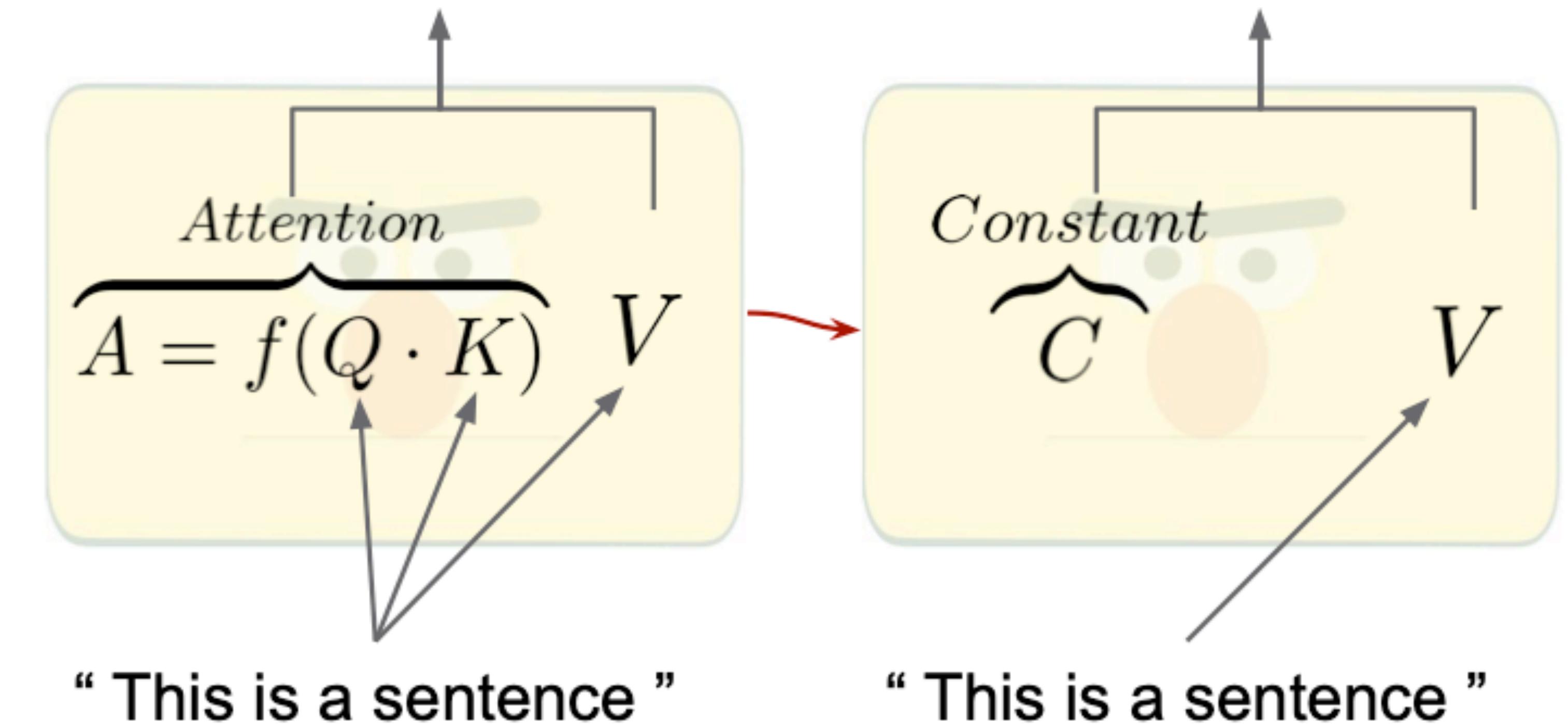




How Much Does Attention Actually Attend?

Hassid, Peng, Rotem, Kasai, Montero, Smith & S., Findings of EMNLP 2022

- Model doesn't collapse
 - Average accuracy loss of 8% only
- Potential for huge savings





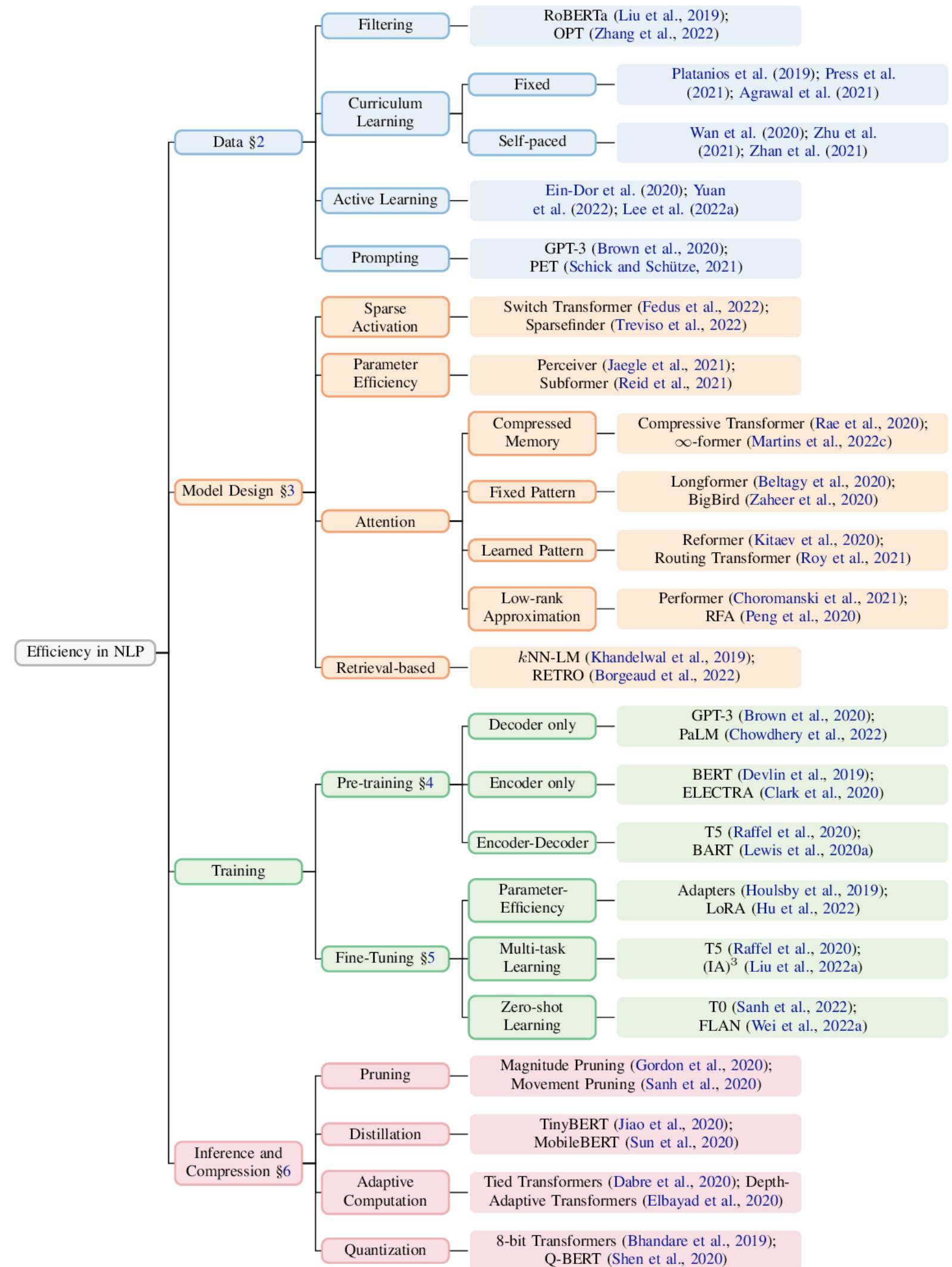
Efficient Modeling

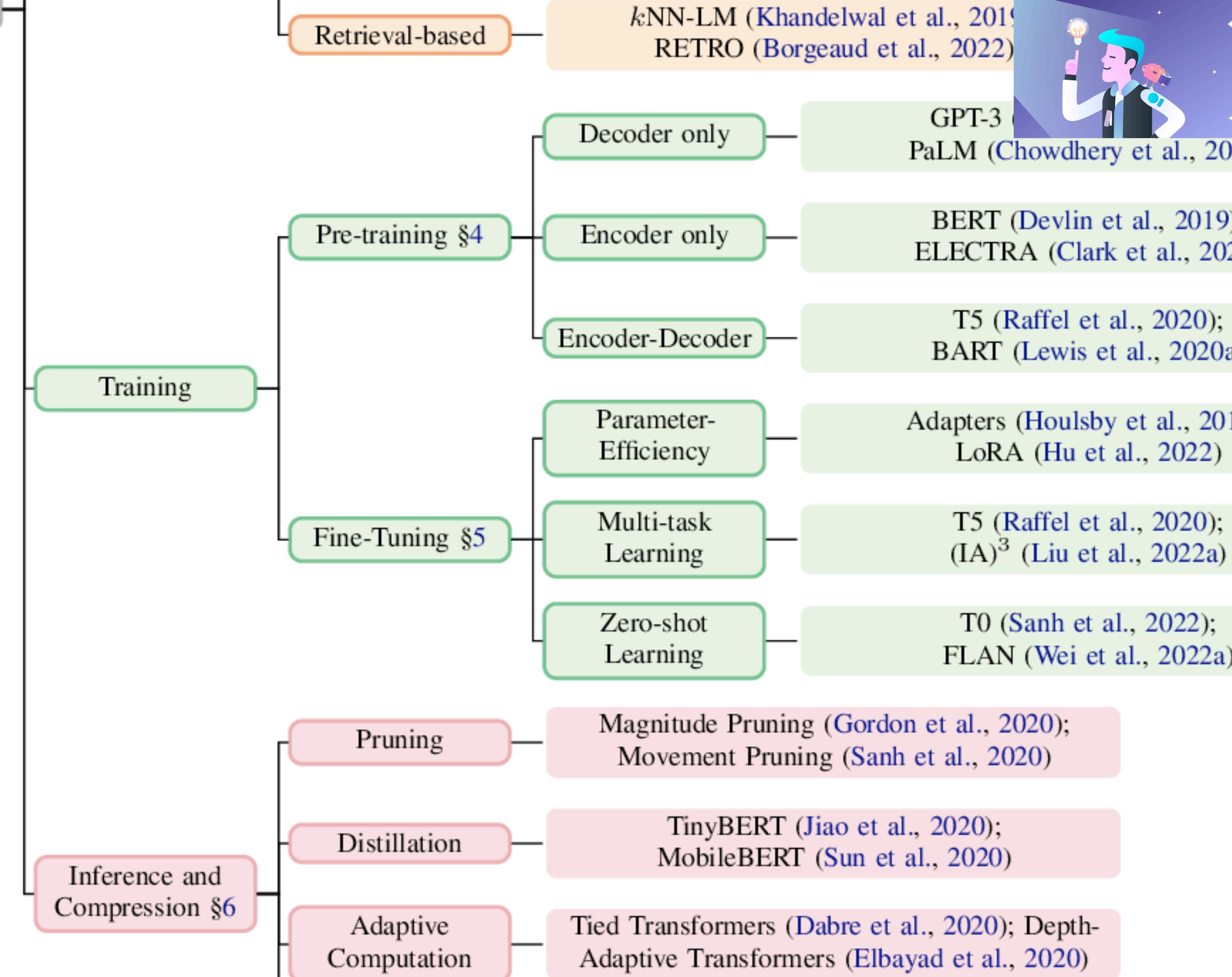
Open Questions

- Can we find the next generation of Transformers?
 - S4 (Gu et al., 2021)
- Should we store knowledge in the model parameters?
 - Retrieval-based models
 - Gu et al (2018); Lewis et al. (2020); Li et al. (2022); Borgeaud et al. (2022)



Efficient Methods in NLP







Space Efficiency

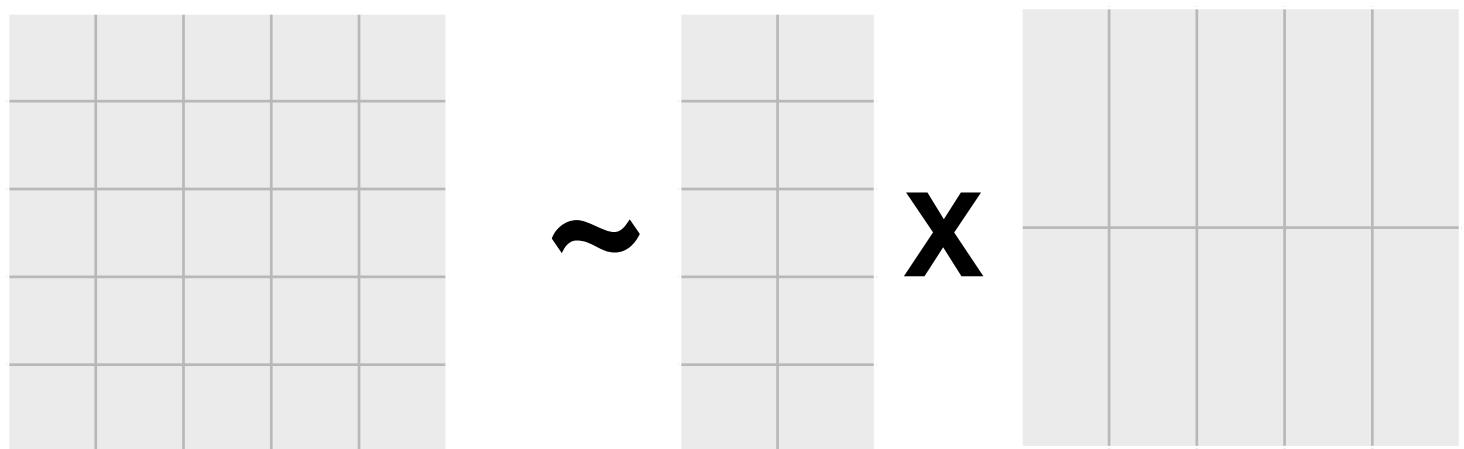
- Weight Factorization
 - Lan et al. (2019); Wang et al. (2019)

$$\begin{matrix} & \sim & \times & \end{matrix}$$

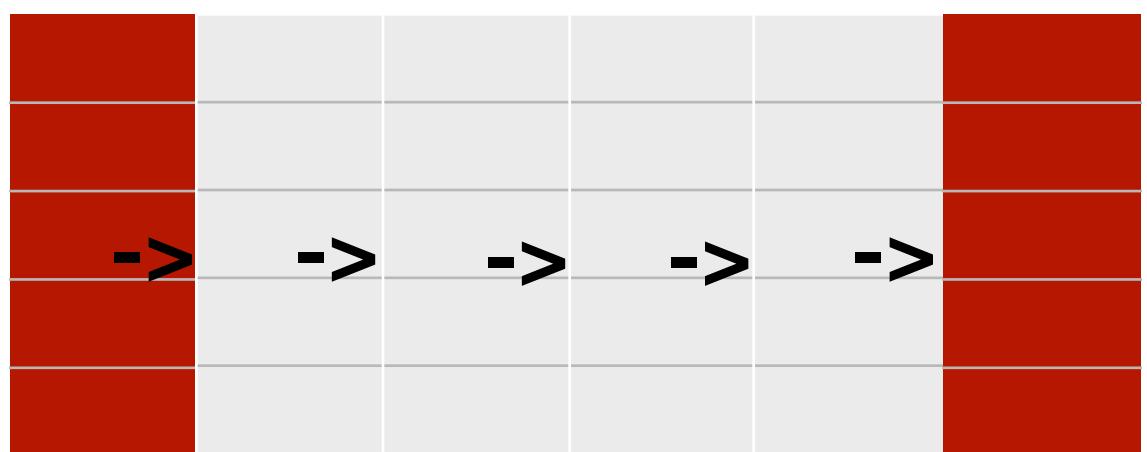


Space Efficiency

- Weight Factorization
 - Lan et al. (2019); Wang et al. (2019)



- Weight Sharing
 - Inan et al. (2016); Press & Wolf (2017); Dehghani et al. (2019)

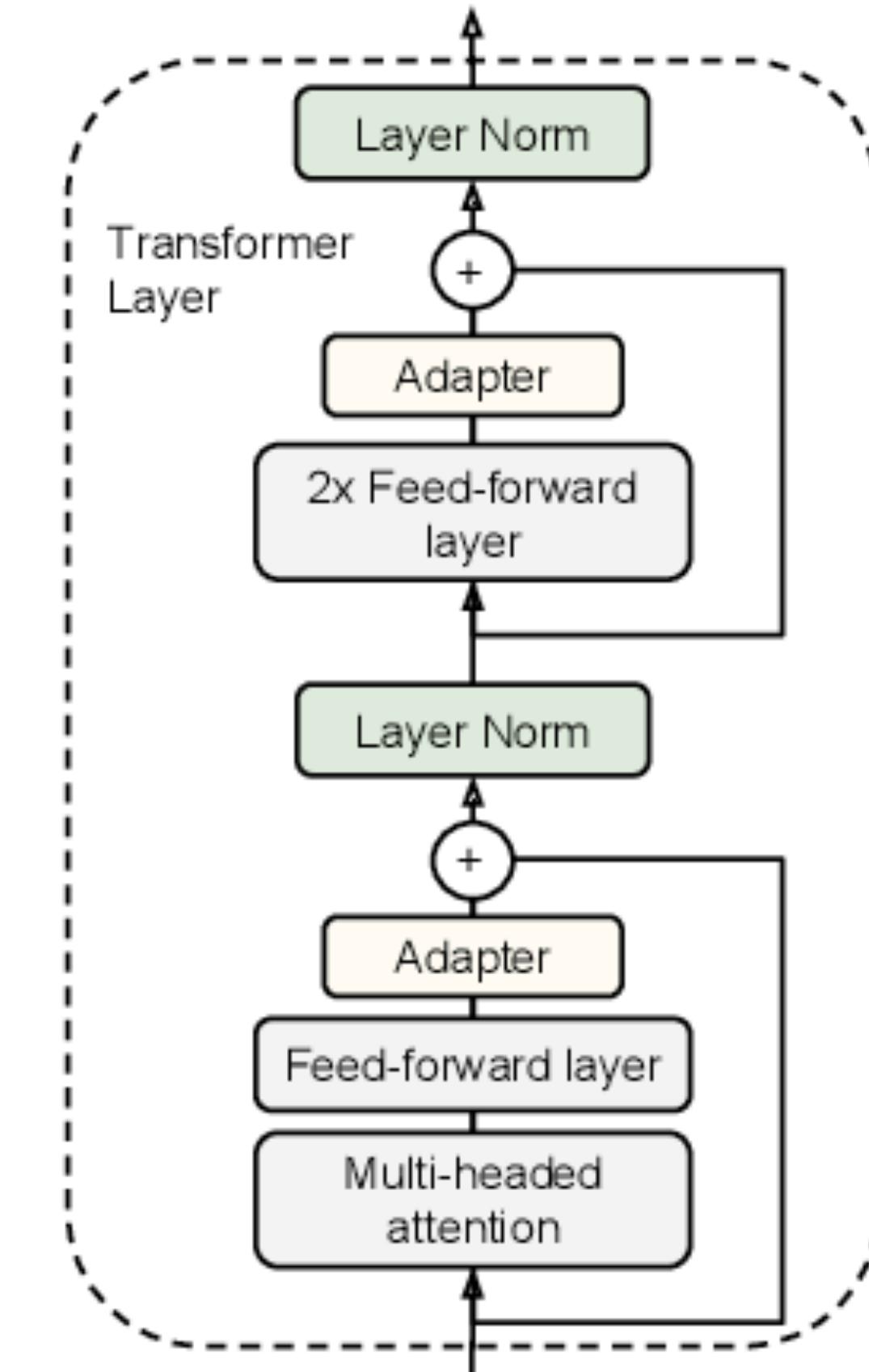




Adapters

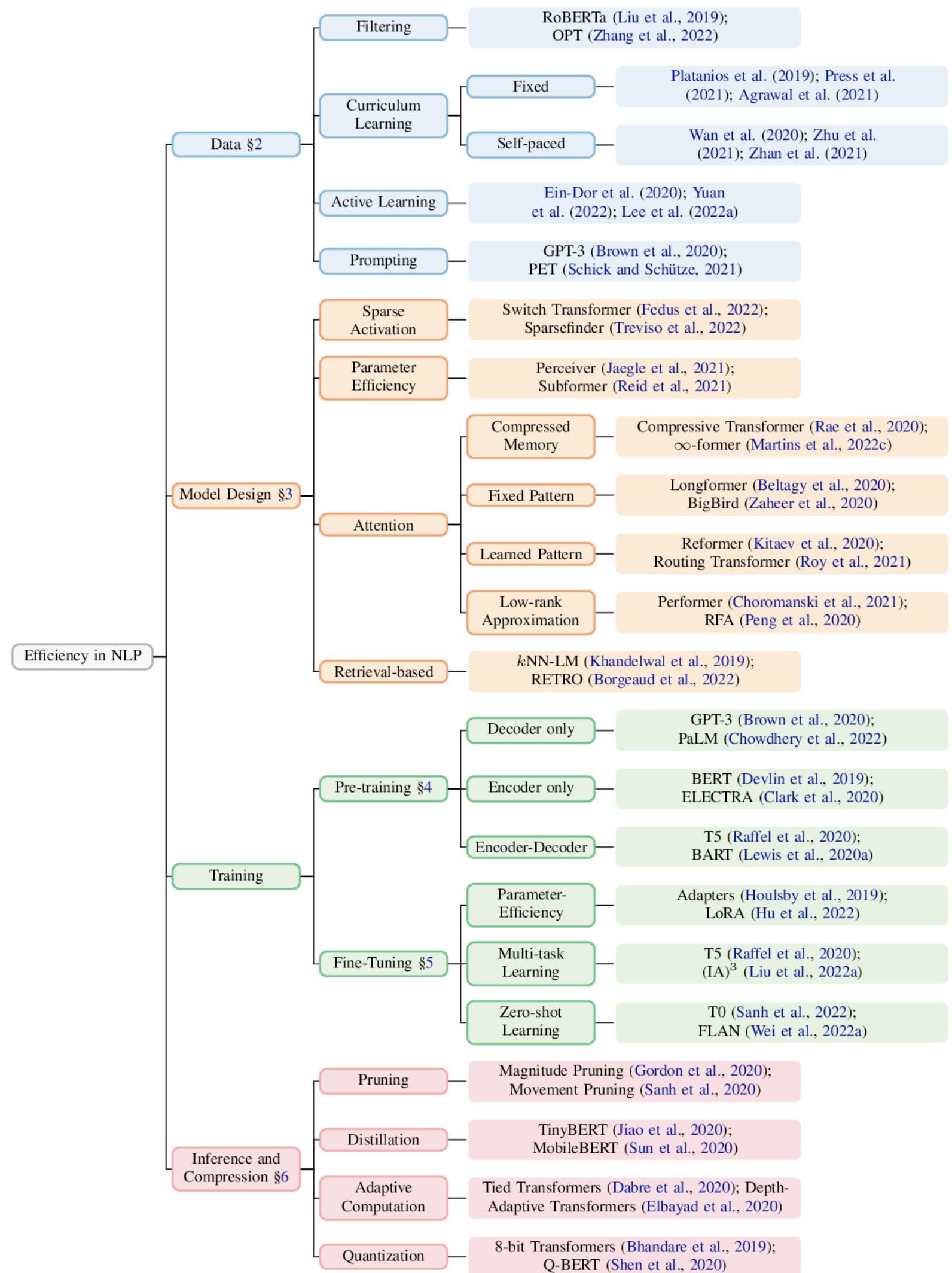
Houlsby et al. (2019)

- Instead of updating all parameters during training, add small components (**adapters**) and only update those





Efficient Methods in NLP





Efficiency

LoR

Fine-Tuning §5

Multi-task Learning

T5 (IA)

Zero-shot Learning

T0 (Sanh et al., 2022);
FLAN (Wei et al., 2022a)

Pruning

Magnitude Pruning (Gordon et al., 2020);
Movement Pruning (Sanh et al., 2020)

Distillation

TinyBERT (Jiao et al., 2020);
MobileBERT (Sun et al., 2020)

Adaptive Computation

Tied Transformers (Dabre et al., 2020); Depth-Adaptive Transformers (Elbayad et al., 2020)

Quantization

8-bit Transformers (Bhandare et al., 2019);
Q-BERT (Shen et al., 2020)

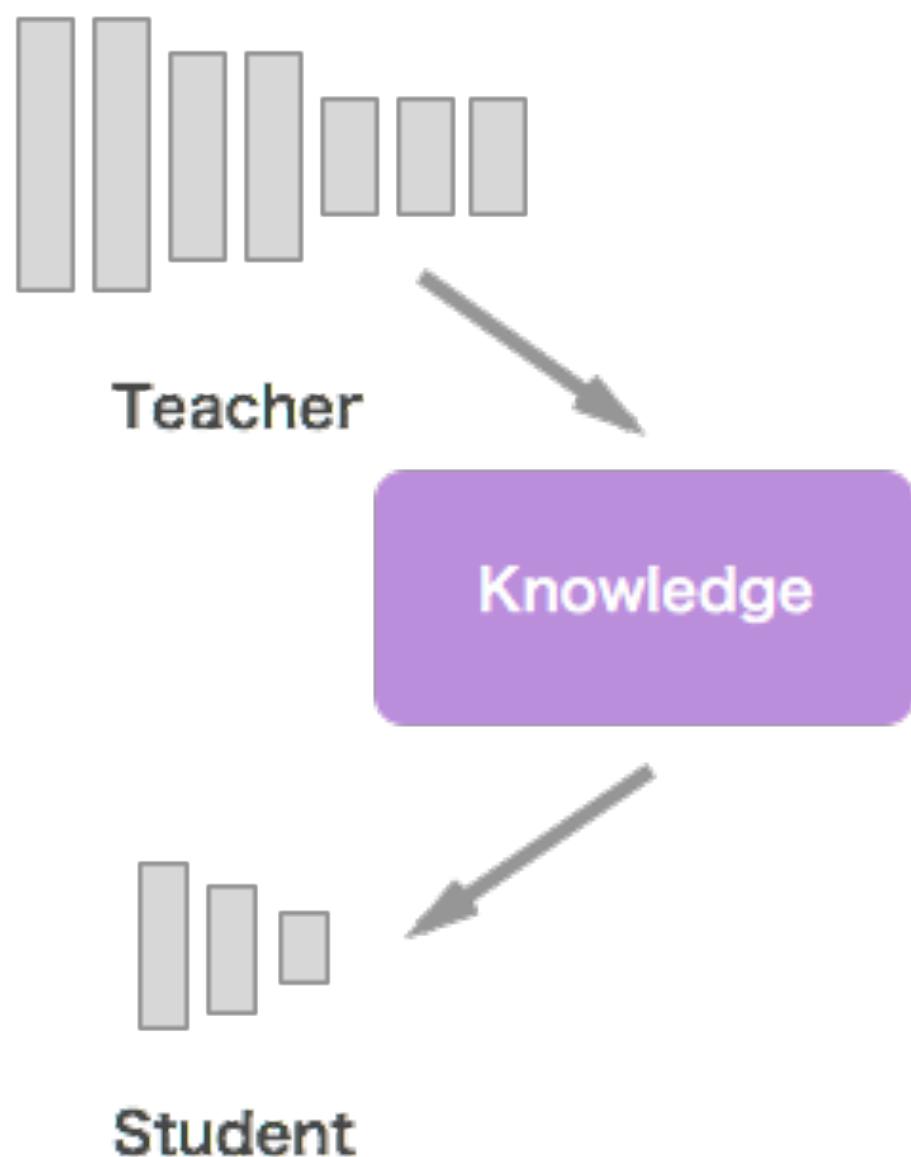
Inference and Compression §6





Model Distillation

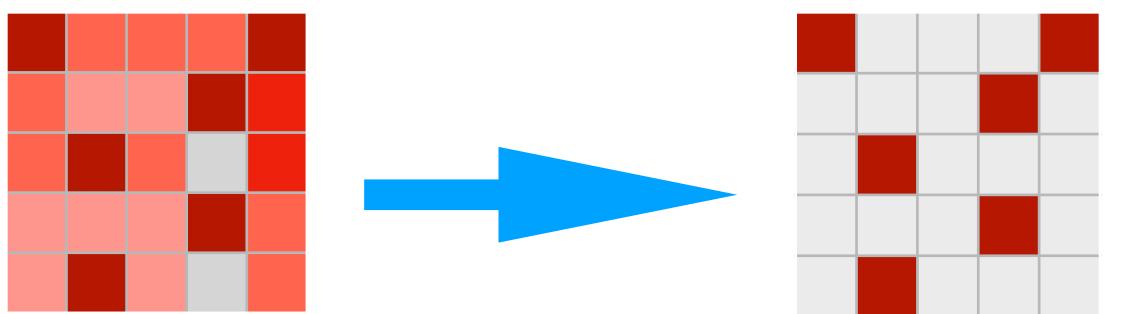
- Student-teacher model
 - Sun et al. (2019); Sanh et al. (2019); Zhao et al (2019)
- Student performs on par with teacher
 - Better than a student of similar size trained from scratch





Pruning

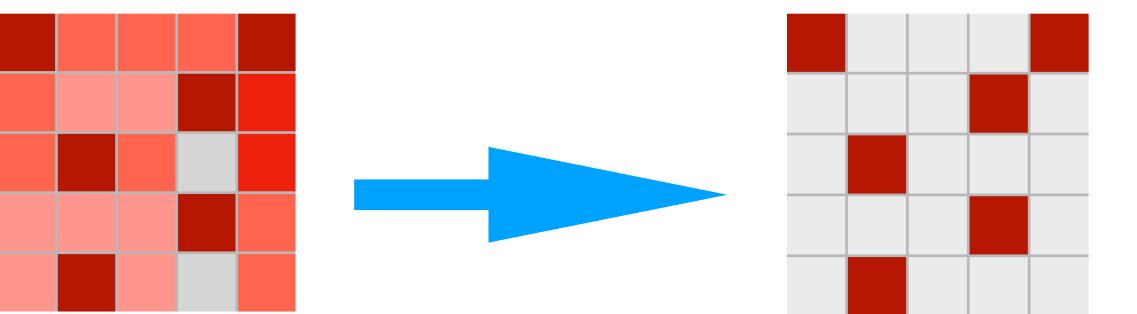
- Learn a sparse version of a trained model
 - Gale et al. (2019); Lee et al. (2019); Frankle & Corbin (2019)





Pruning

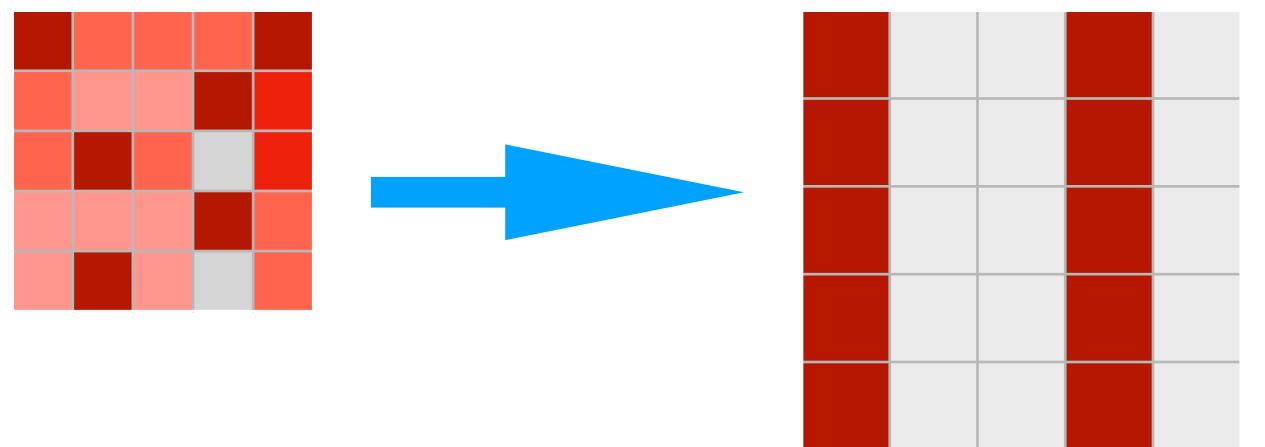
- Learn a sparse version of a trained model
 - Gale et al. (2019); Lee et al. (2019); Frankle & Corbin (2019)
- Hard to translate to speed and energy savings
 - Although see BlockSparse (Gray et al., 2017)
 - Potentially useful for fast training (Dettmers & Zettlemoyer, 2019)





Structured Pruning

- Learn a sparse structure of the trained model
 - I.e., prune attention heads, layers, etc.
 - Gordon et al. (2018); Dodge, S. et al. (2019); Michel et al. (2019); Wang et al. (2019); Fan et al. (2020)
- Could be leveraged to **perform fewer computations**





Quantization

- Reduced precision
 - Zafrir et al. (2019); Shen et al. (2019); Bhandare et al. (2019)

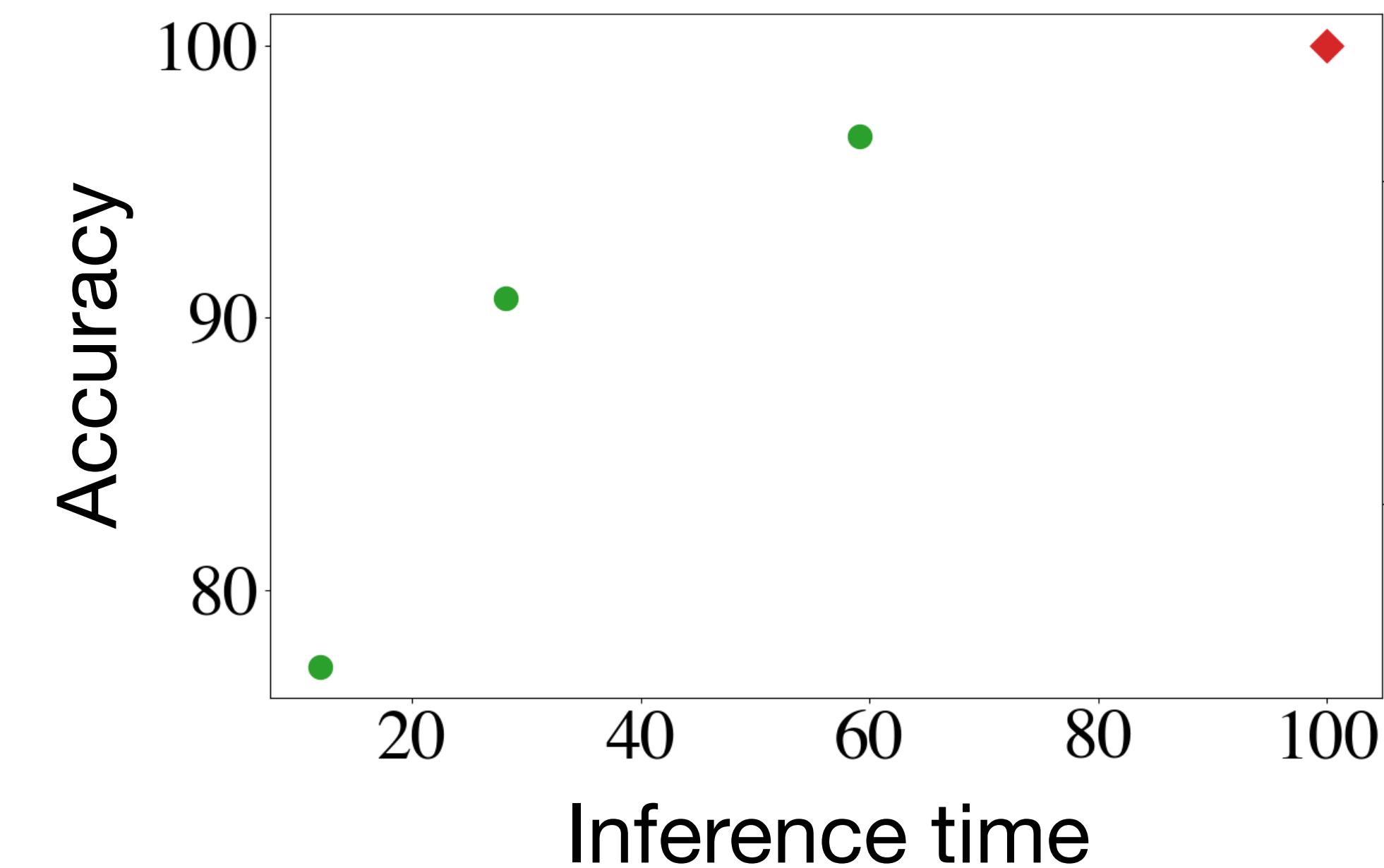




Matching Model and Instance Complexity

S. et al., ACL 2020

*Run an **efficient** model on “easy” instances,
and an **expensive** model on “hard” instances*

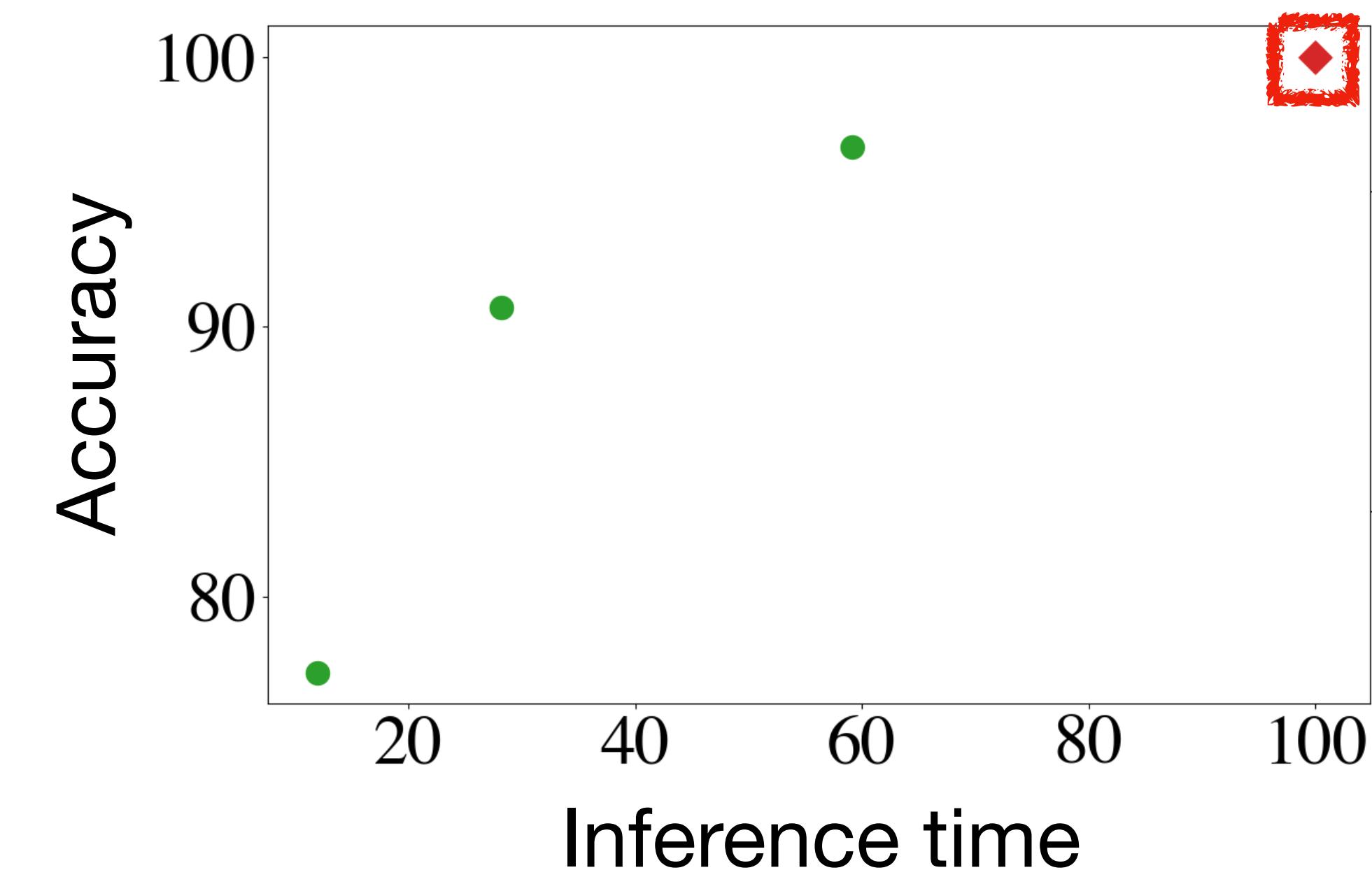




Matching Model and Instance Complexity

S. et al., ACL 2020

*Run an **efficient** model on “easy” instances,
and an **expensive** model on “hard” instances*

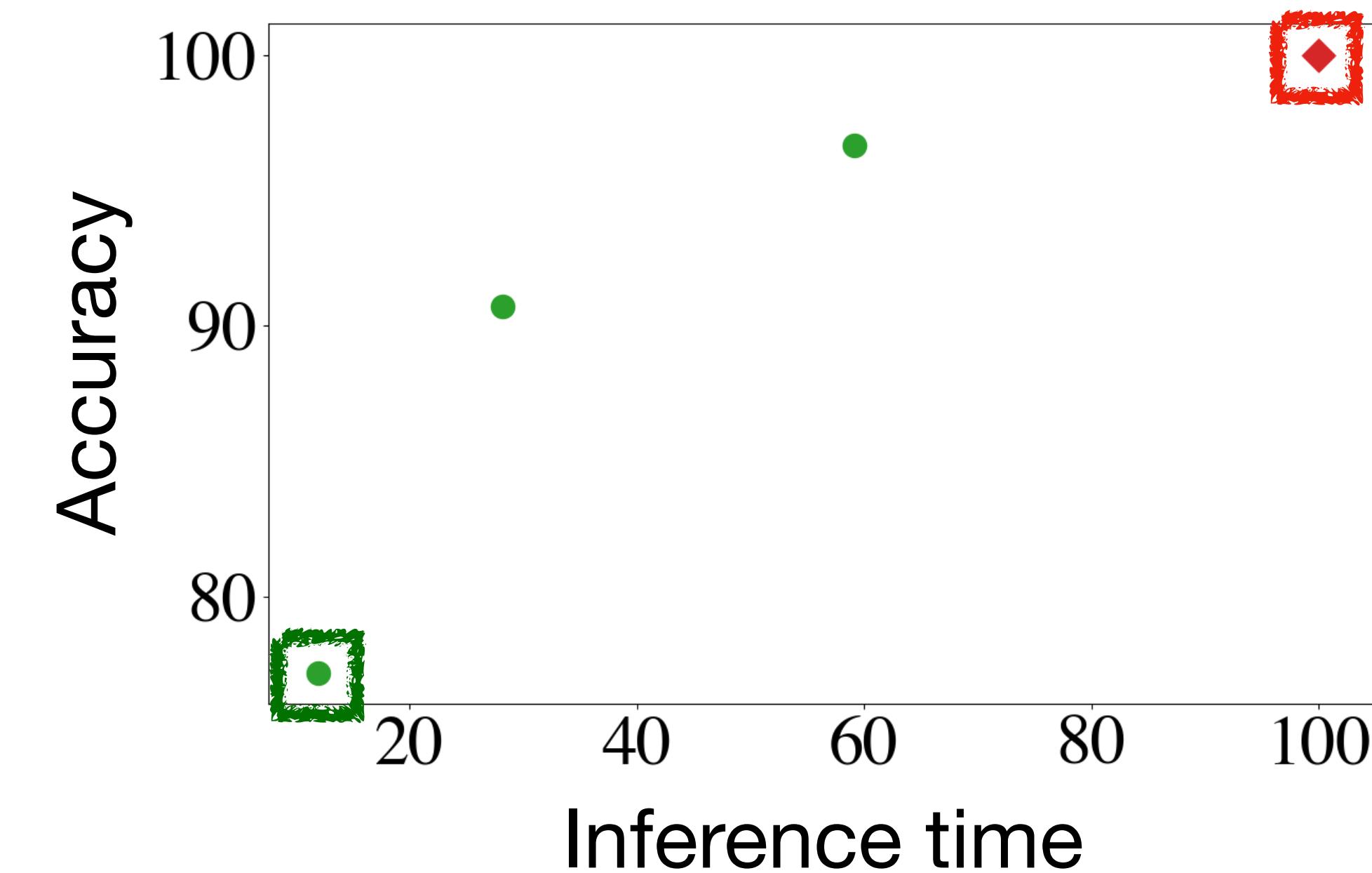




Matching Model and Instance Complexity

S. et al., ACL 2020

*Run an **efficient** model on “easy” instances,
and an **expensive** model on “hard” instances*

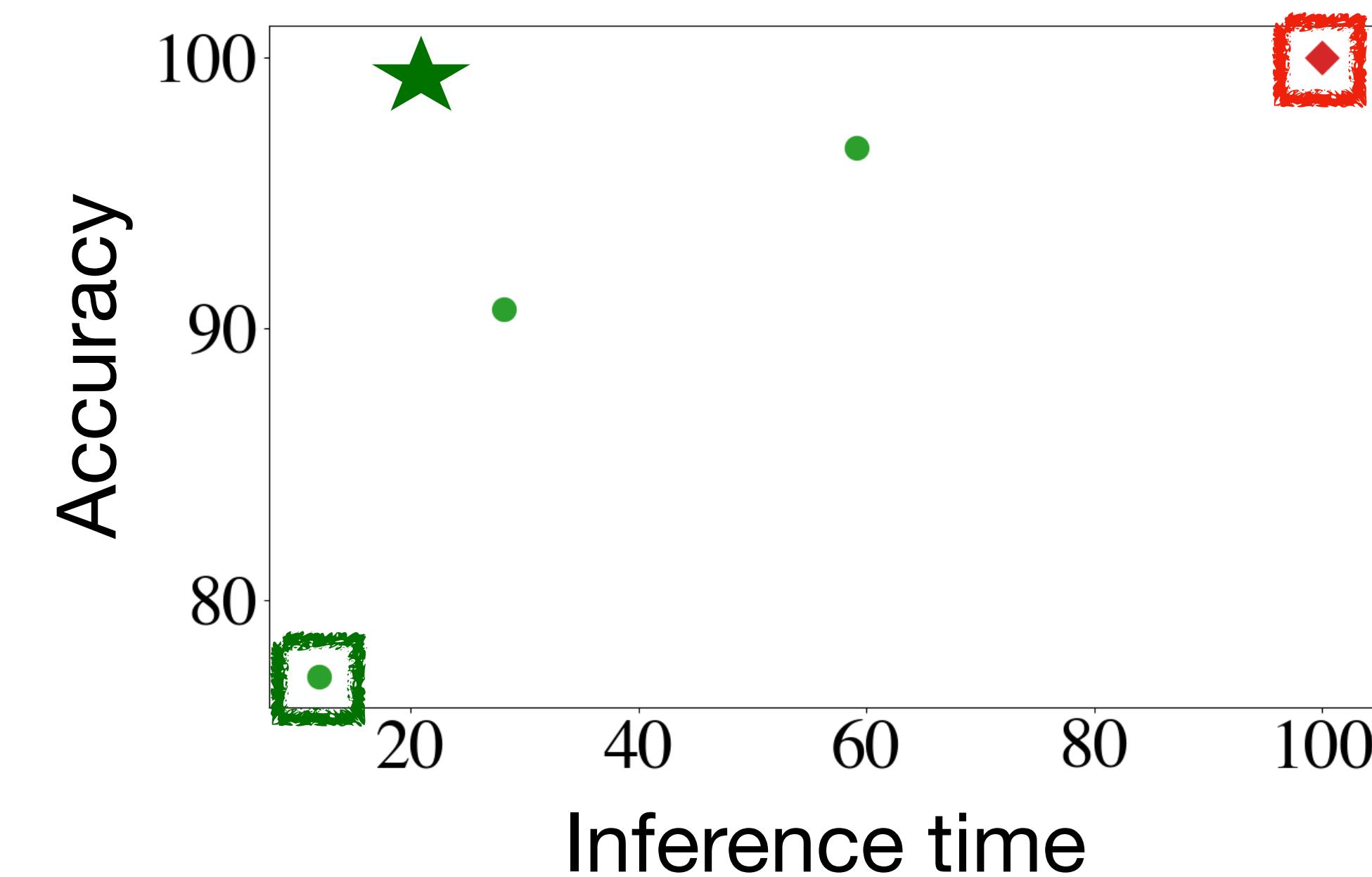




Matching Model and Instance Complexity

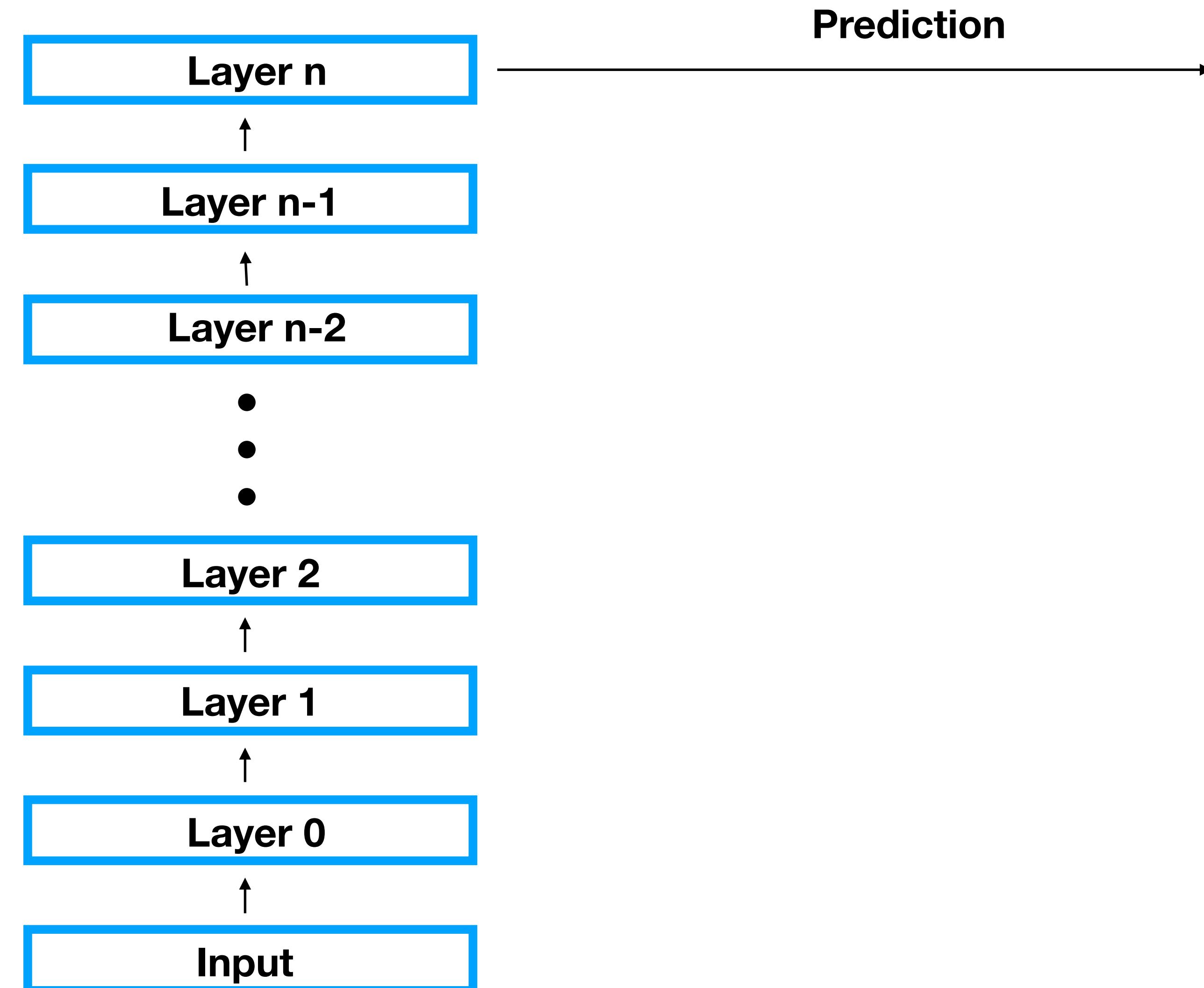
S. et al., ACL 2020

*Run an **efficient** model on “easy” instances,
and an **expensive** model on “hard” instances*



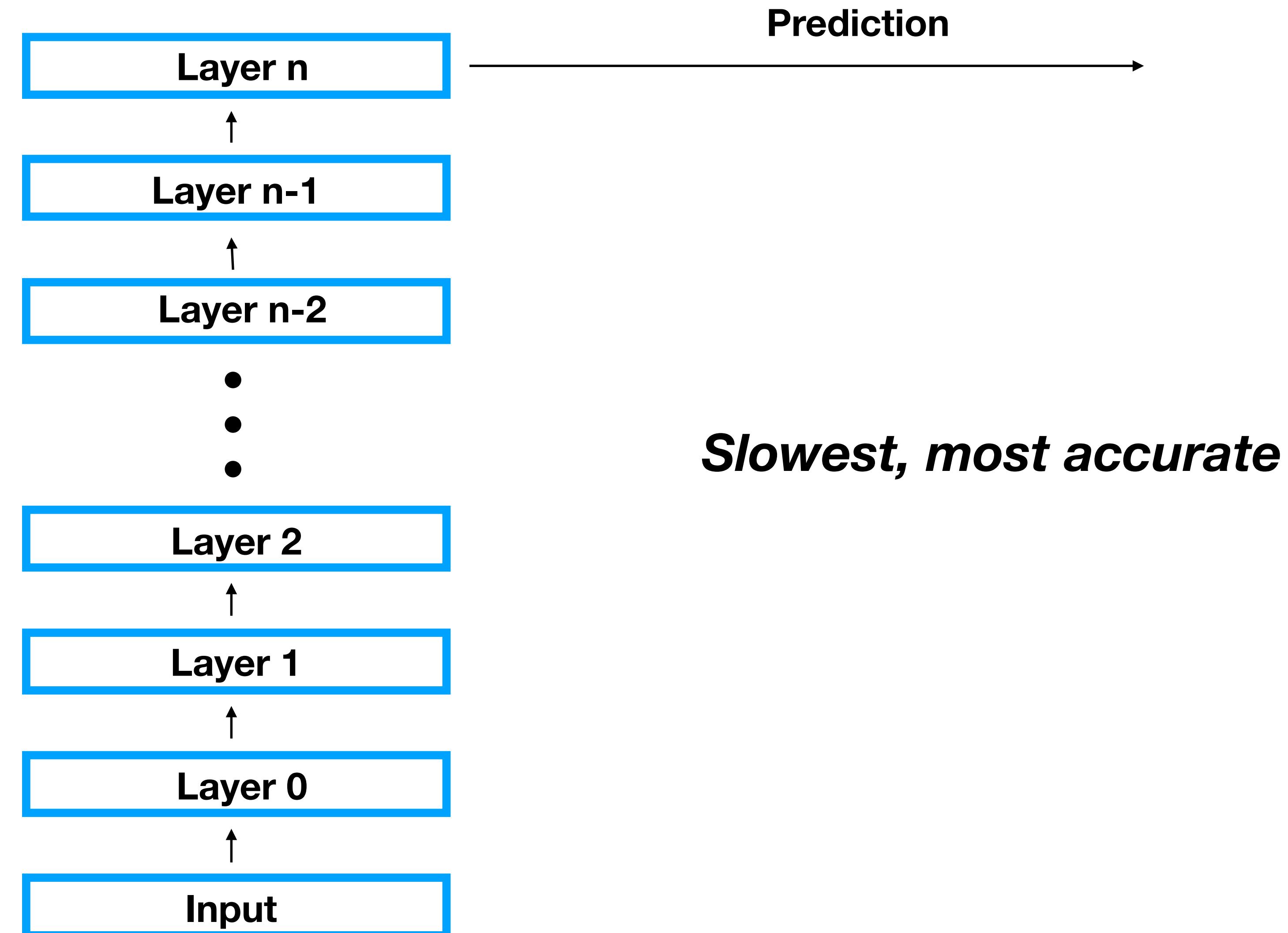


Pretrained BERT Fine-tuning



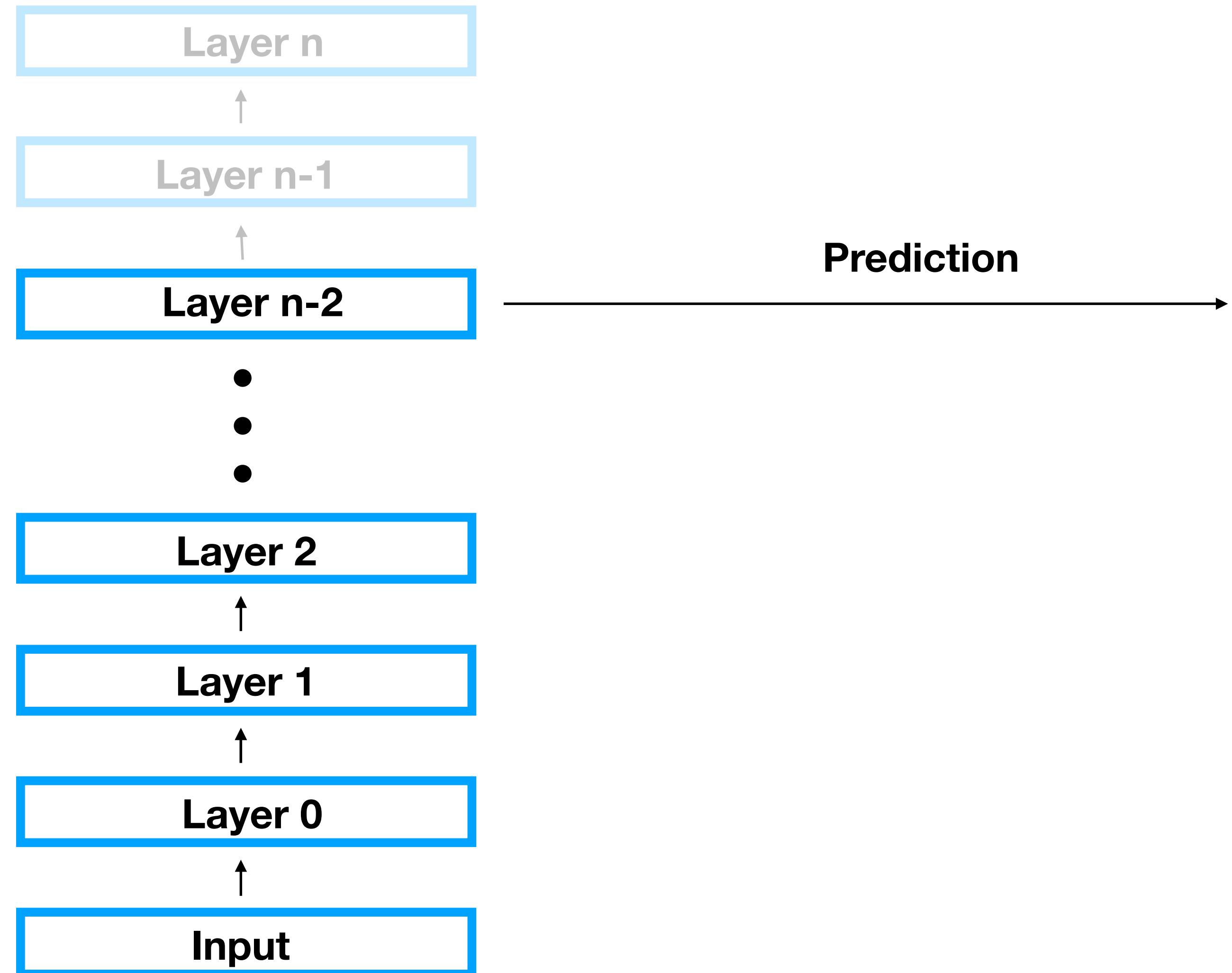


Pretrained BERT Fine-tuning



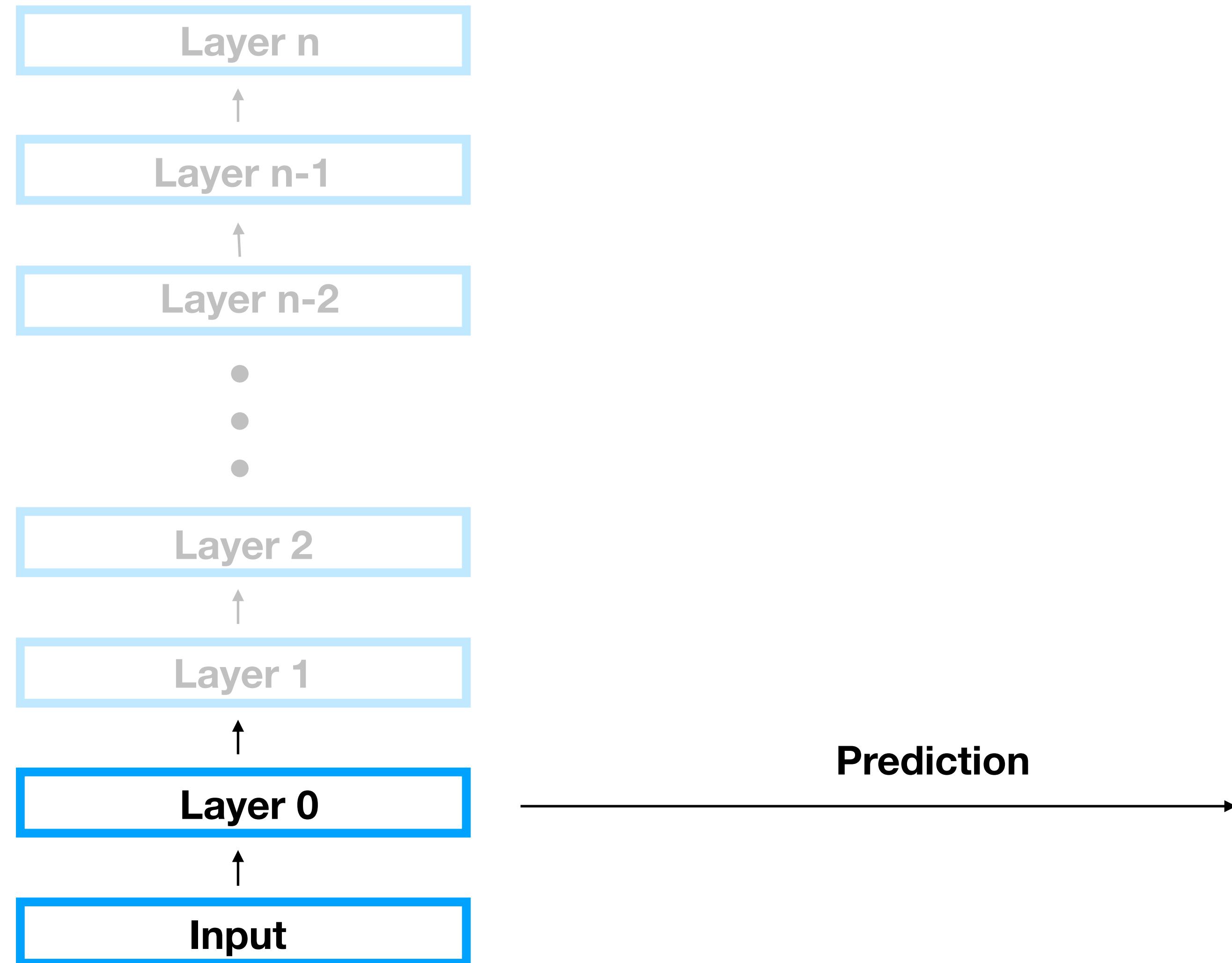


Faster, less Accurate



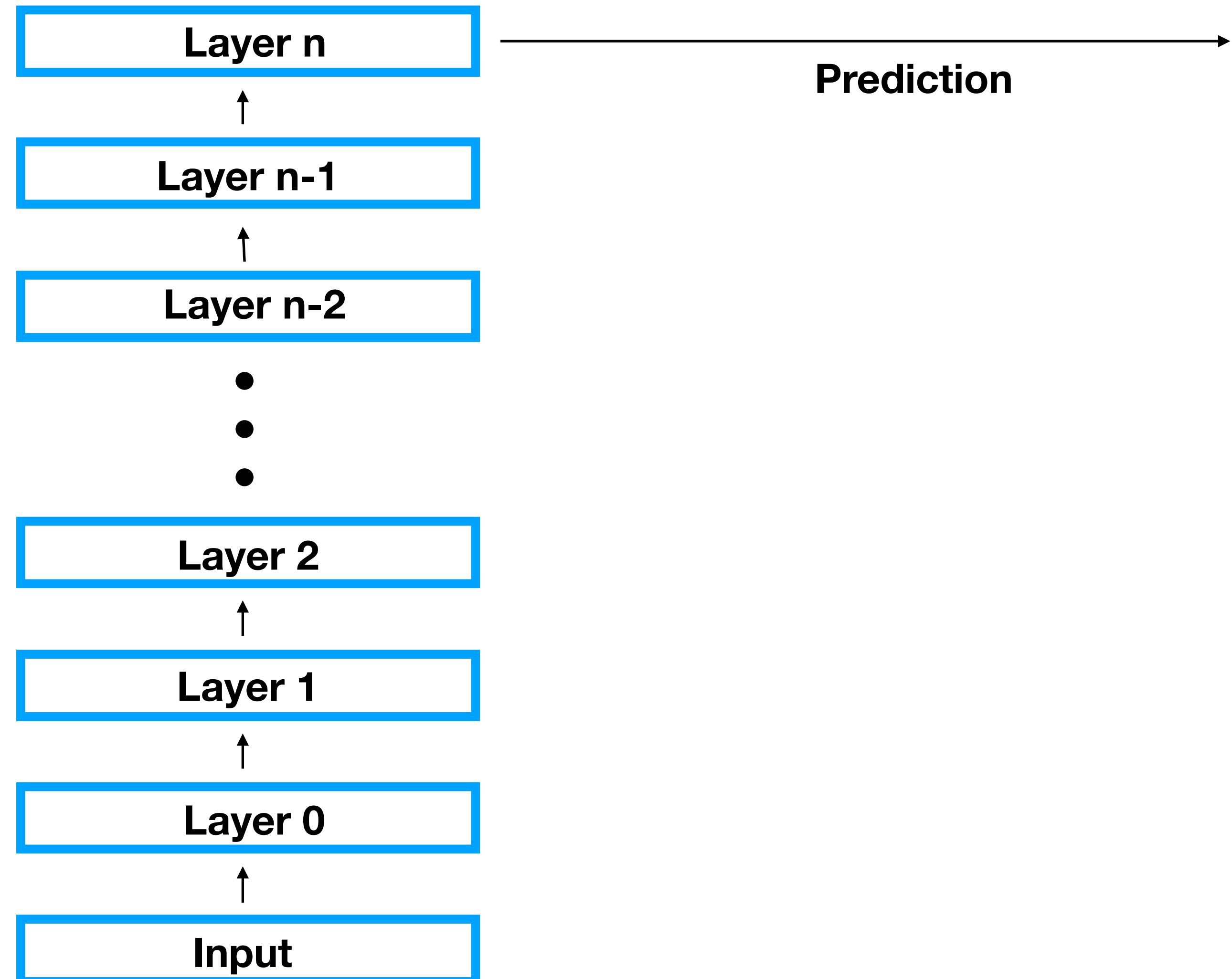


Fastest, least Accurate



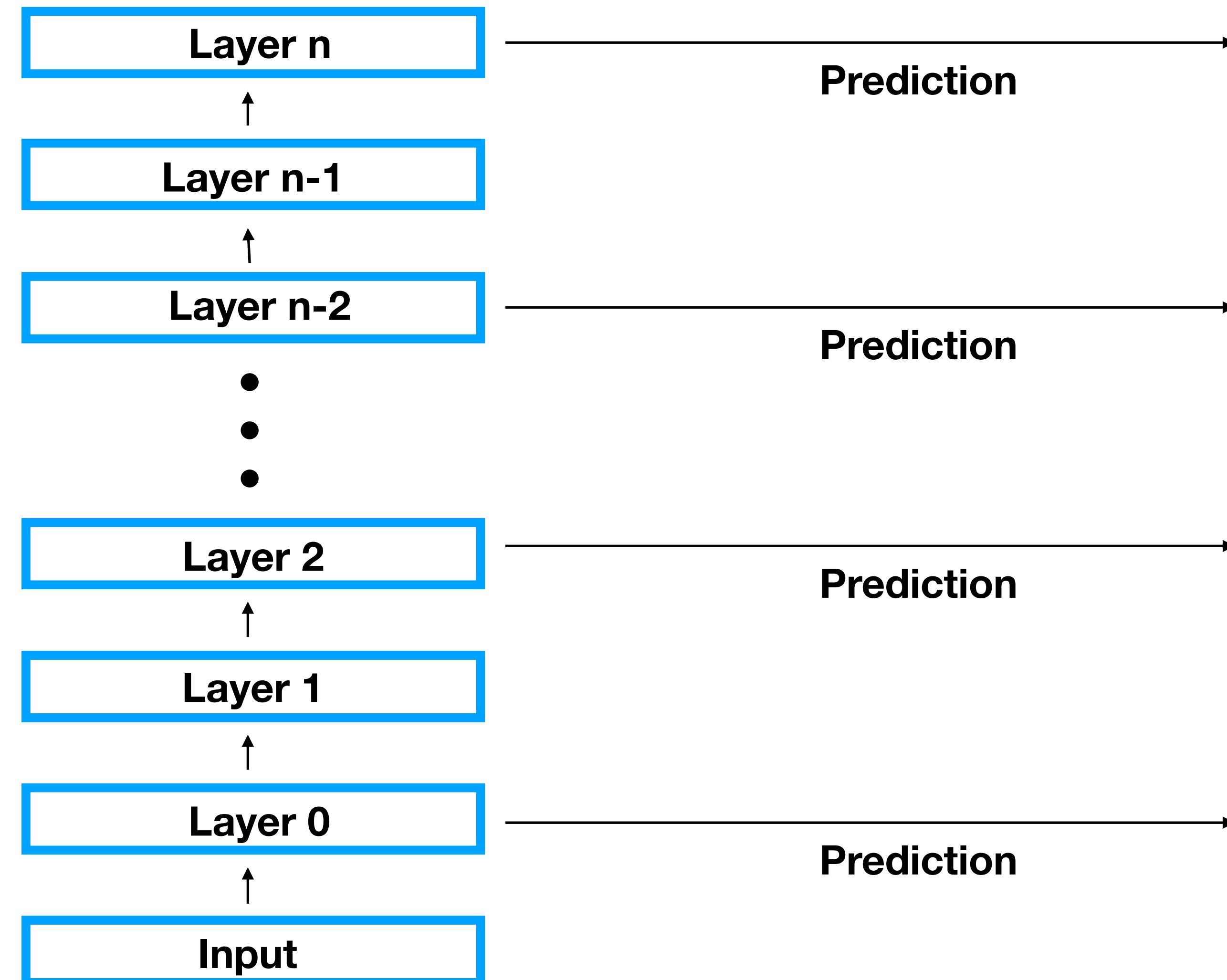


Our Approach: Training Time



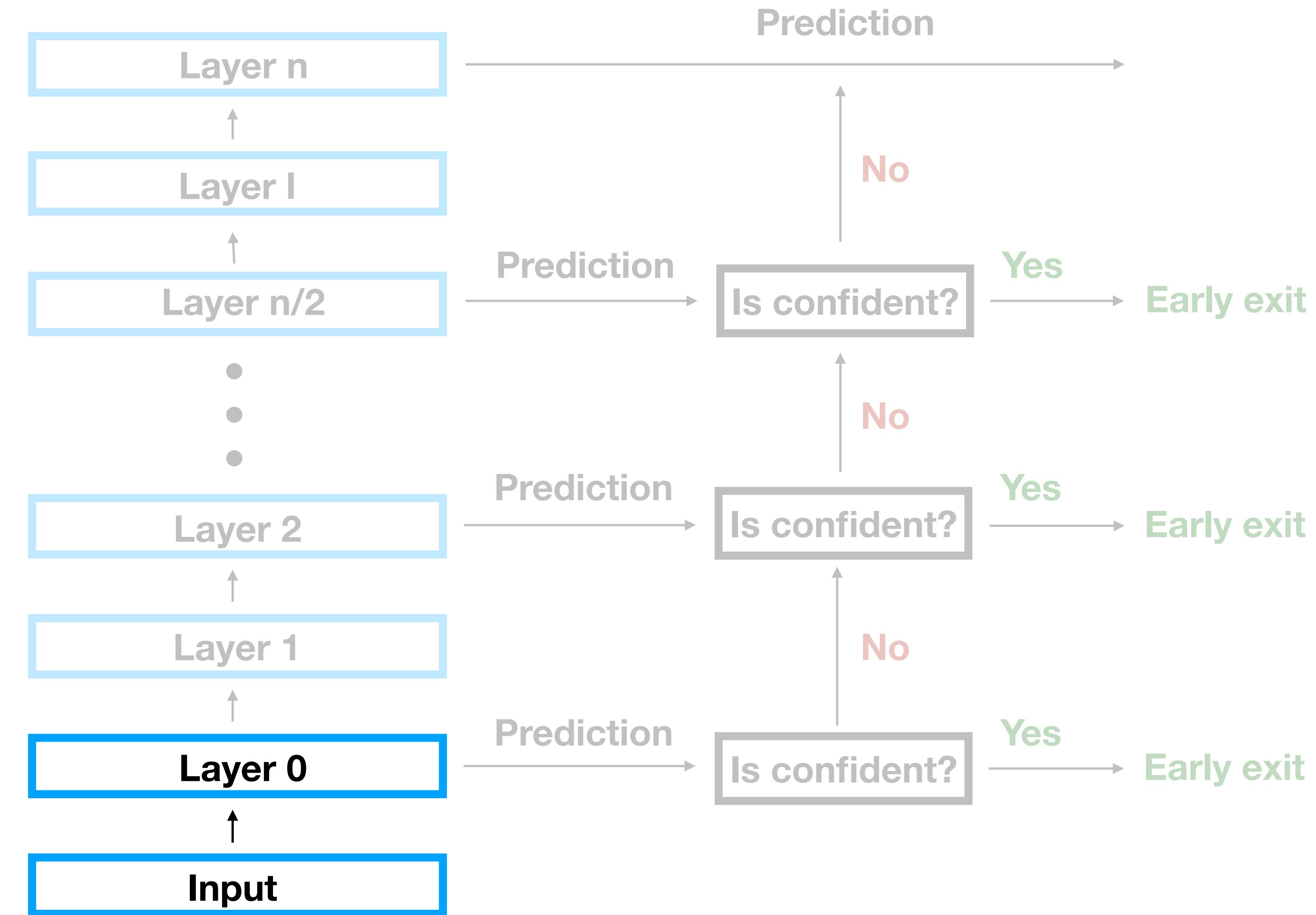


Our Approach: Training Time



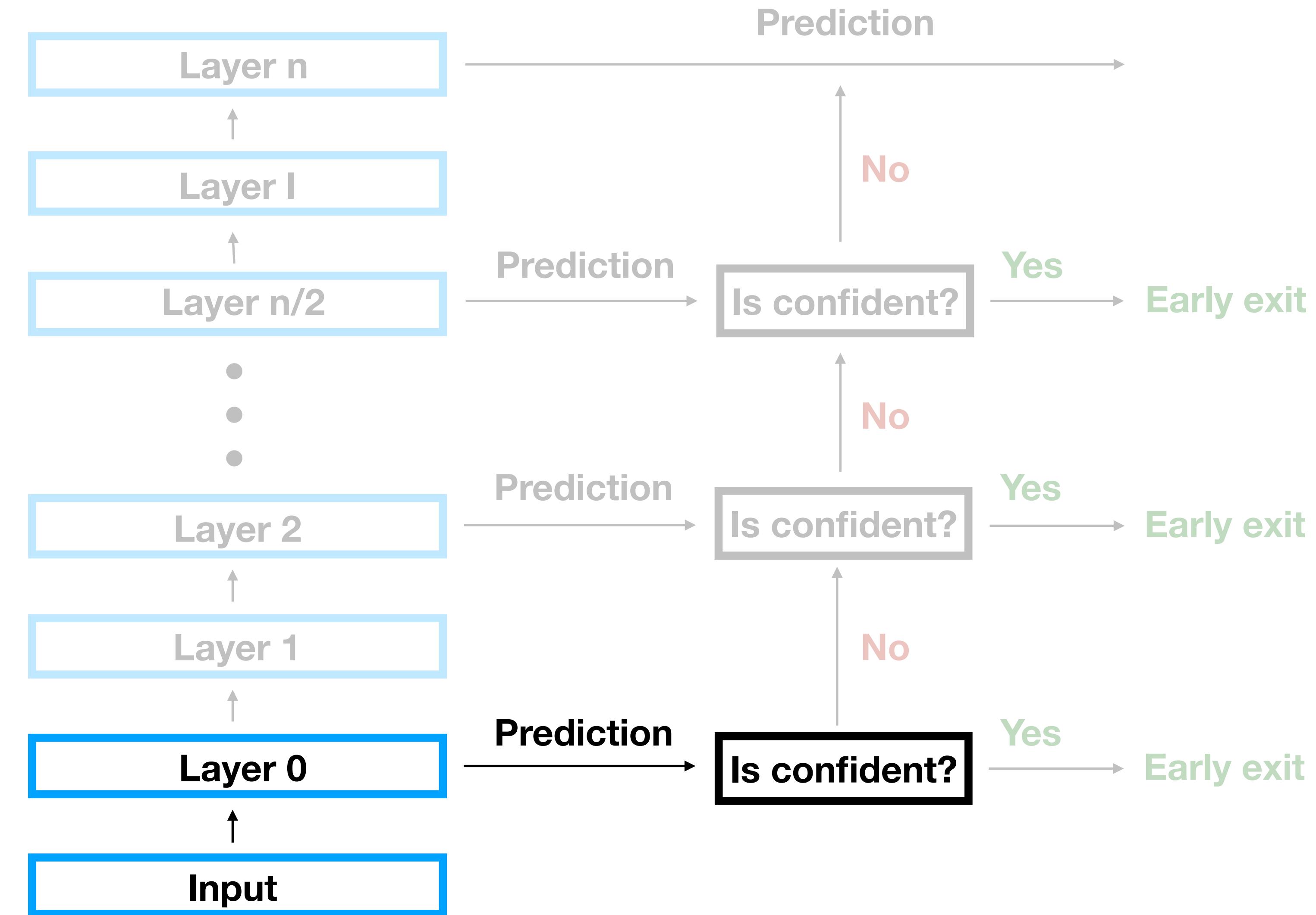


Our Approach: Test Time



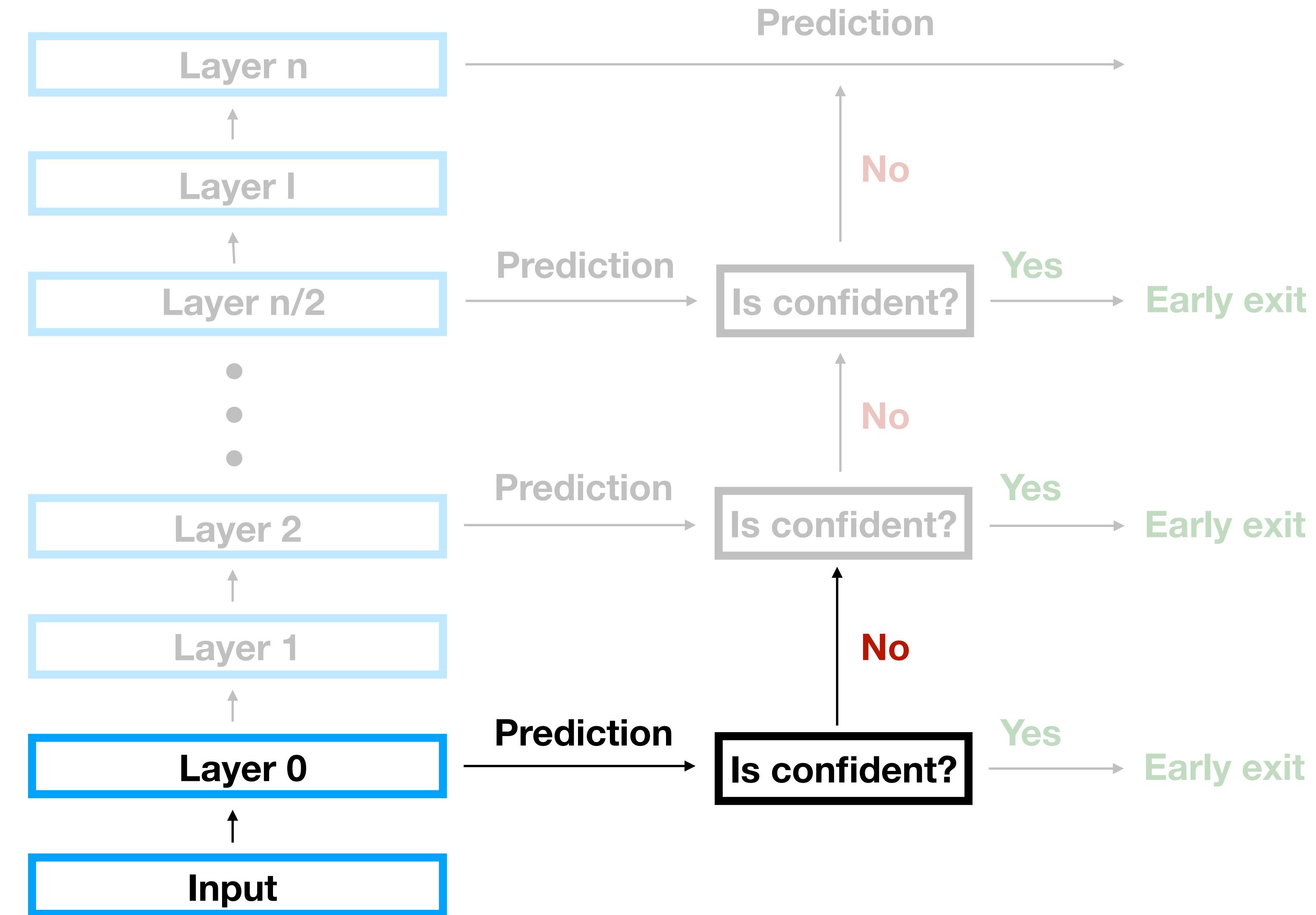


Our Approach: Test Time



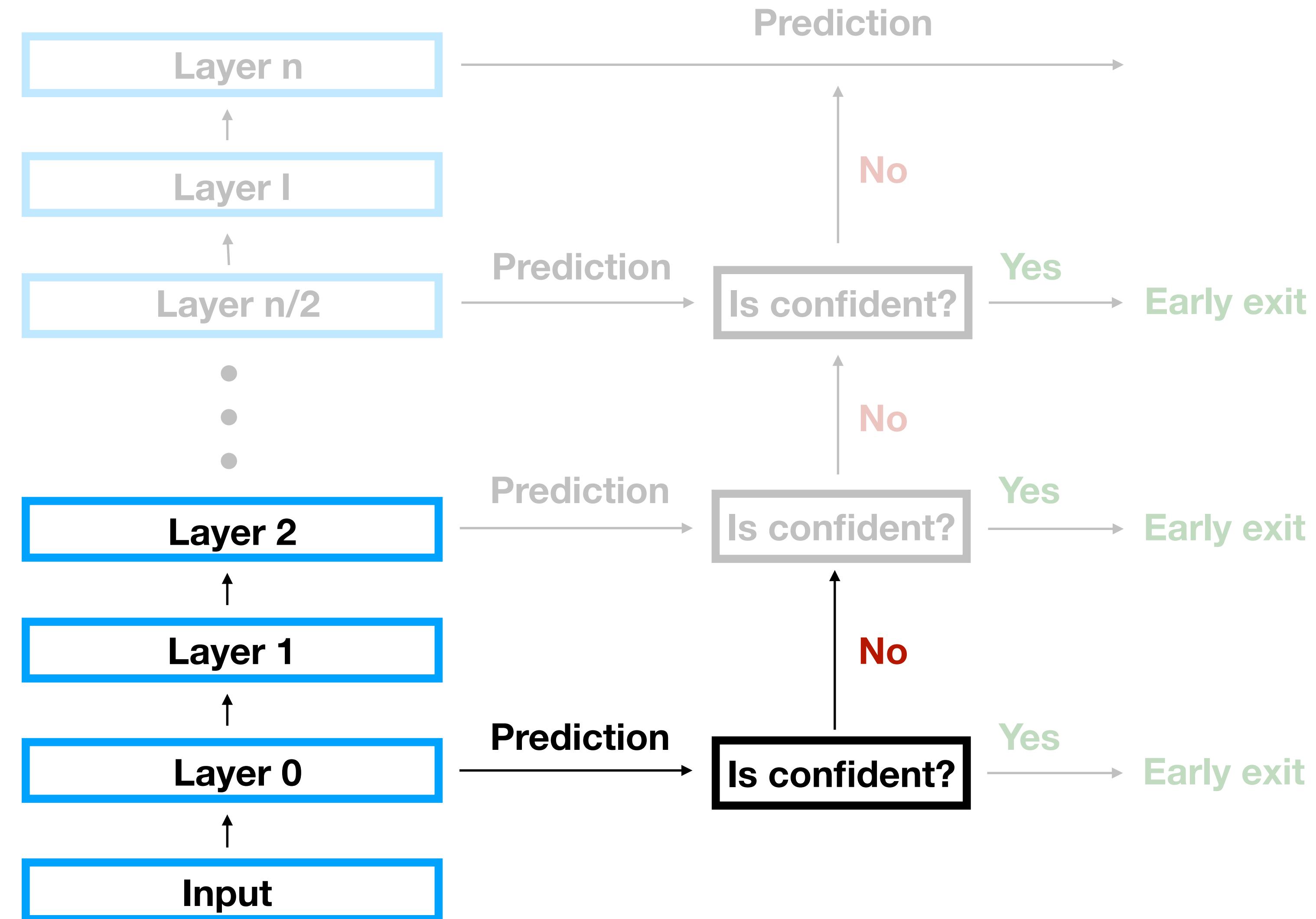


Our Approach: Test Time



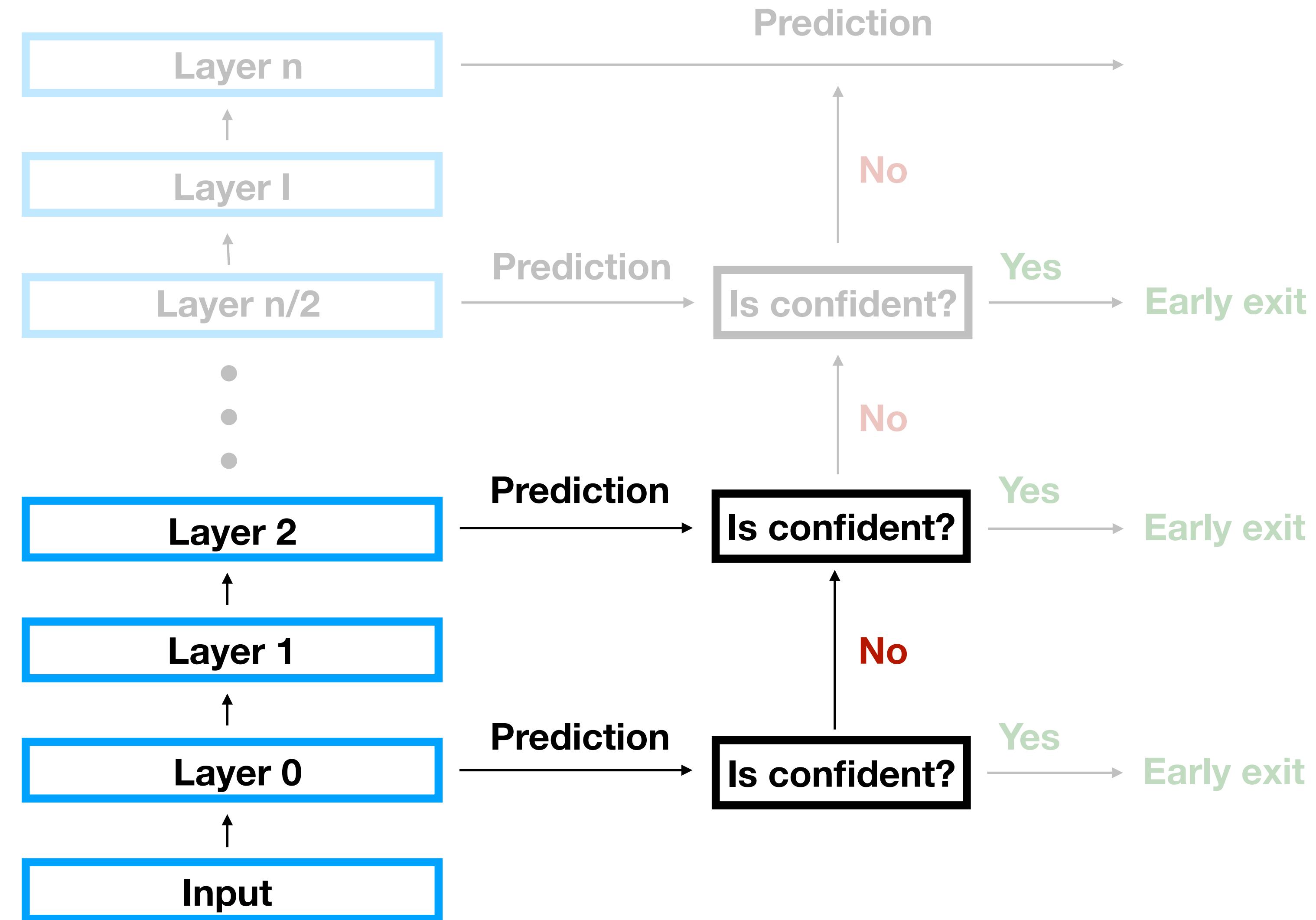


Our Approach: Test Time



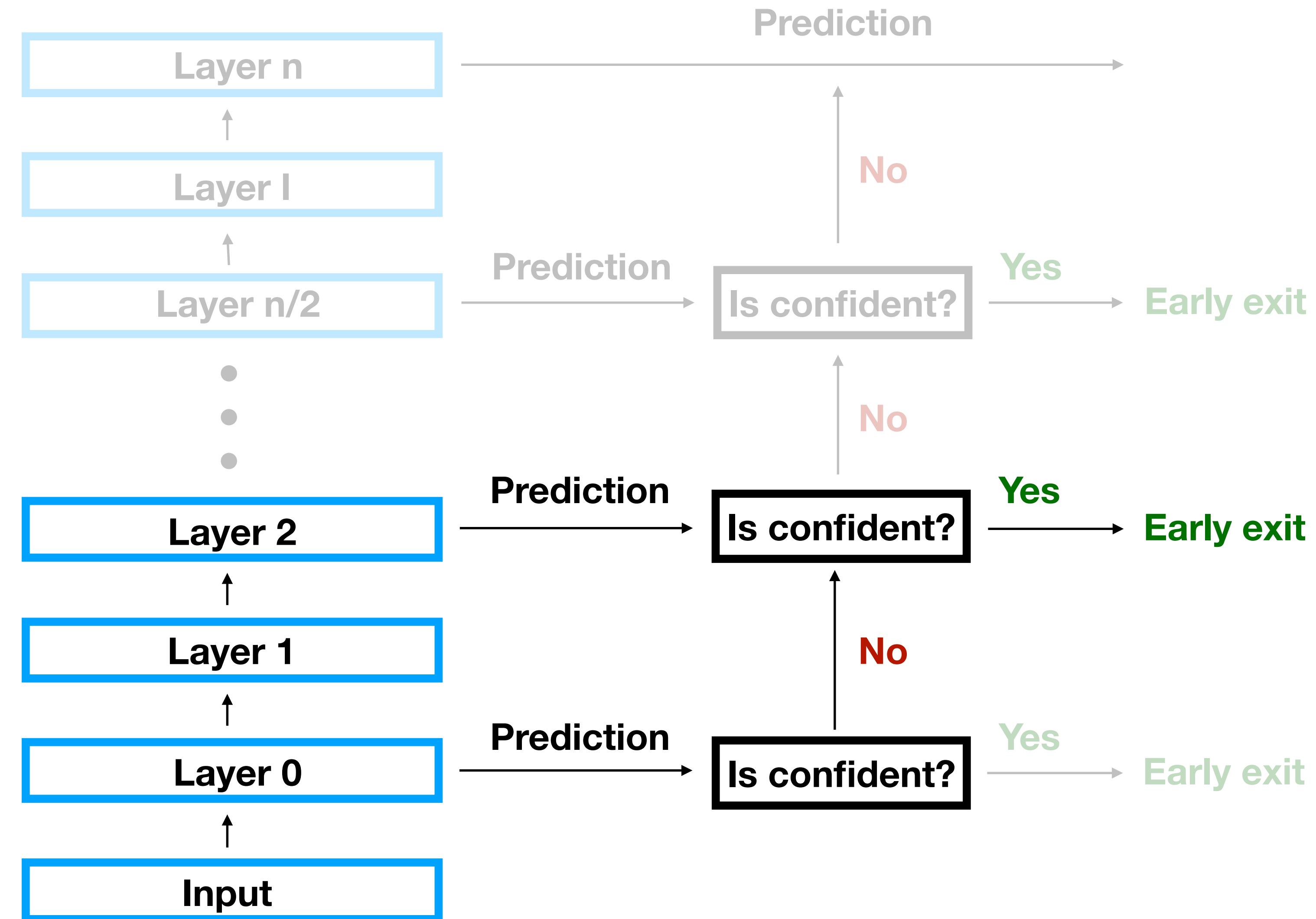


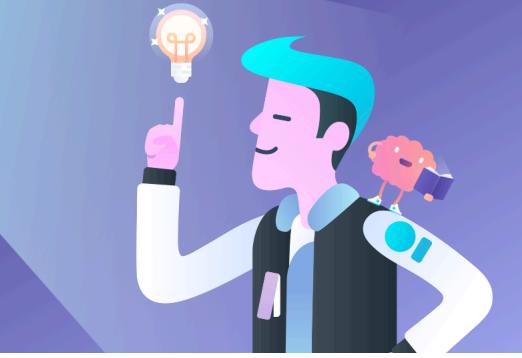
Our Approach: Test Time



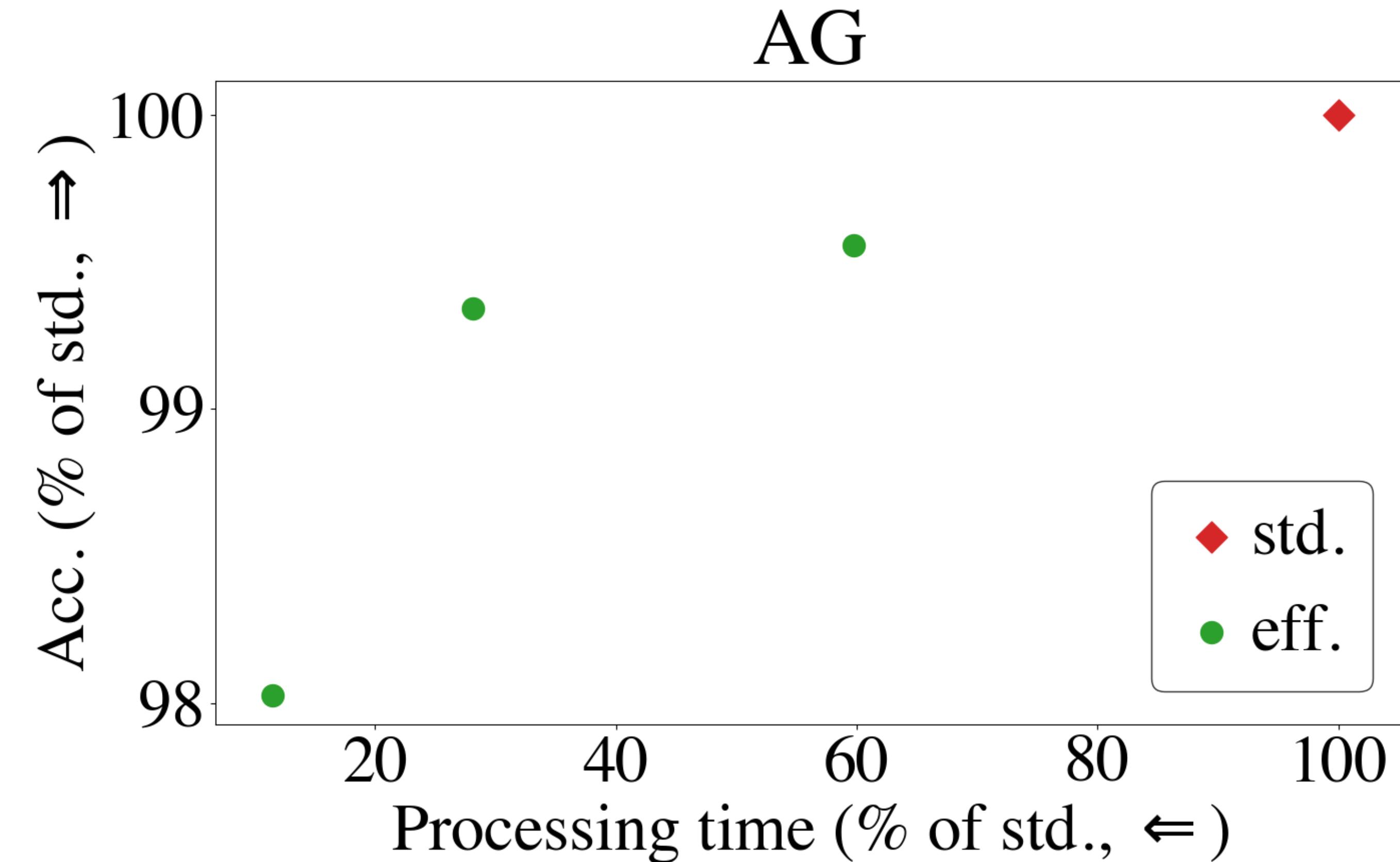


Our Approach: Test Time



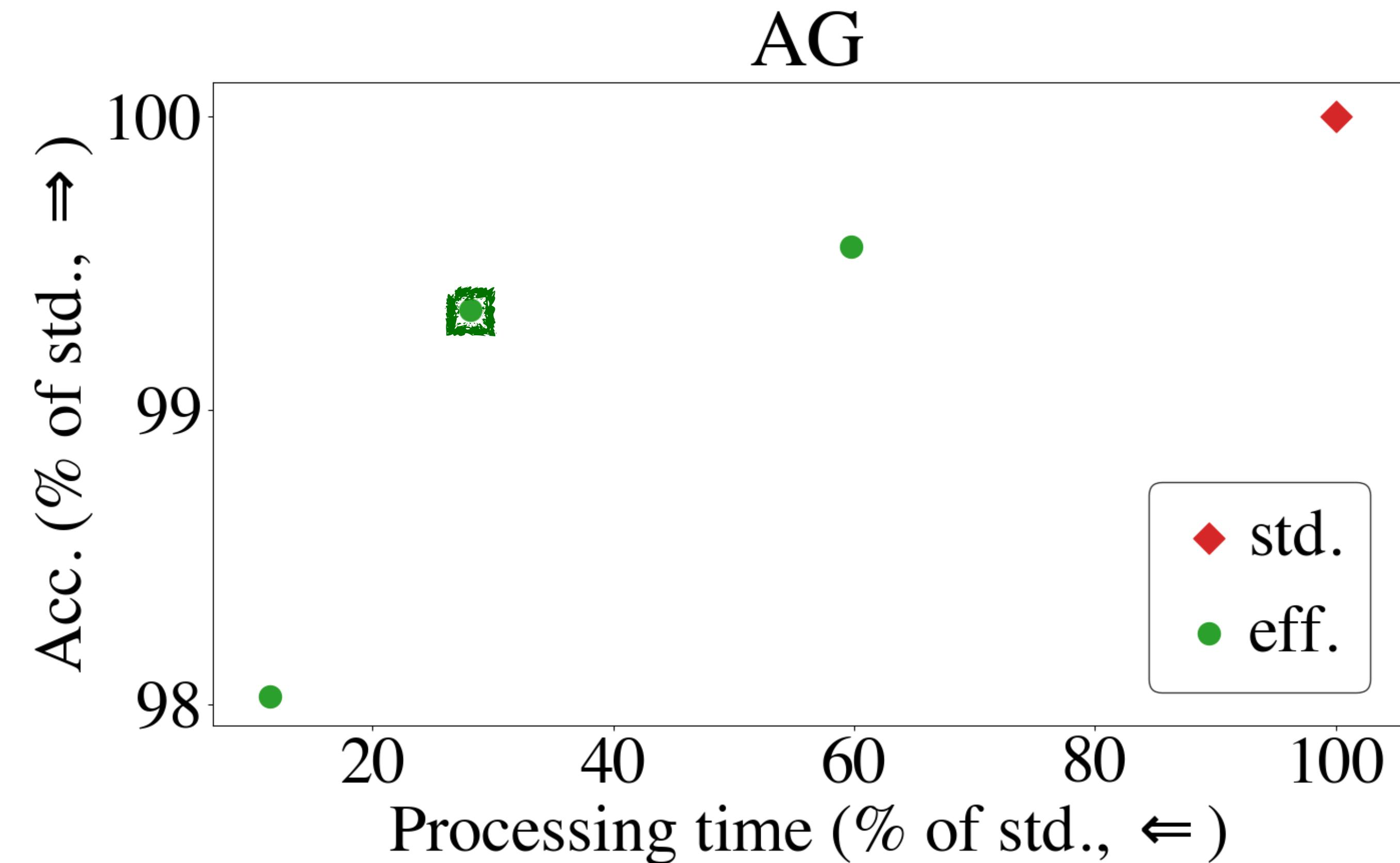


Strong Baselines!



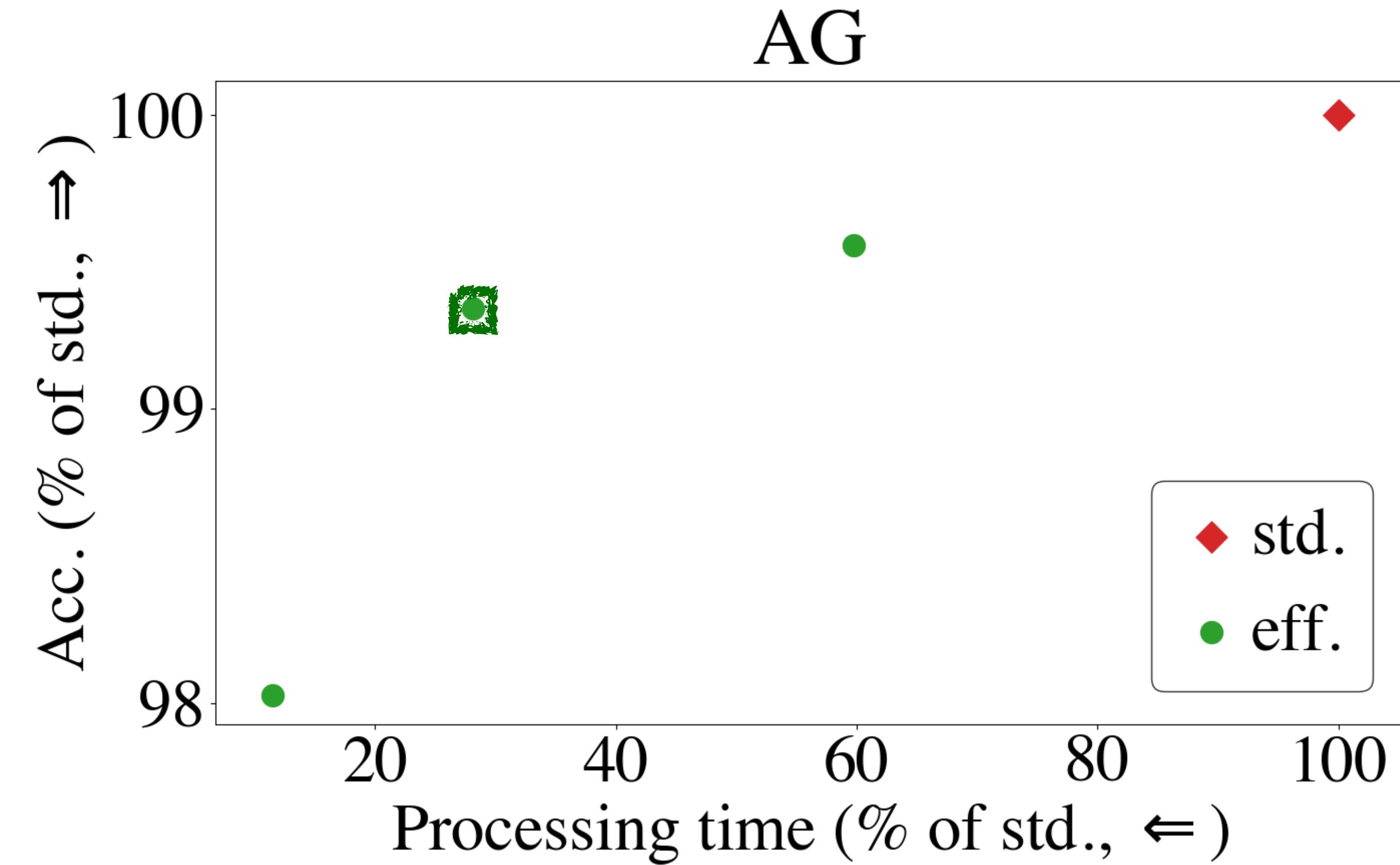


Strong Baselines!





Strong Baselines!

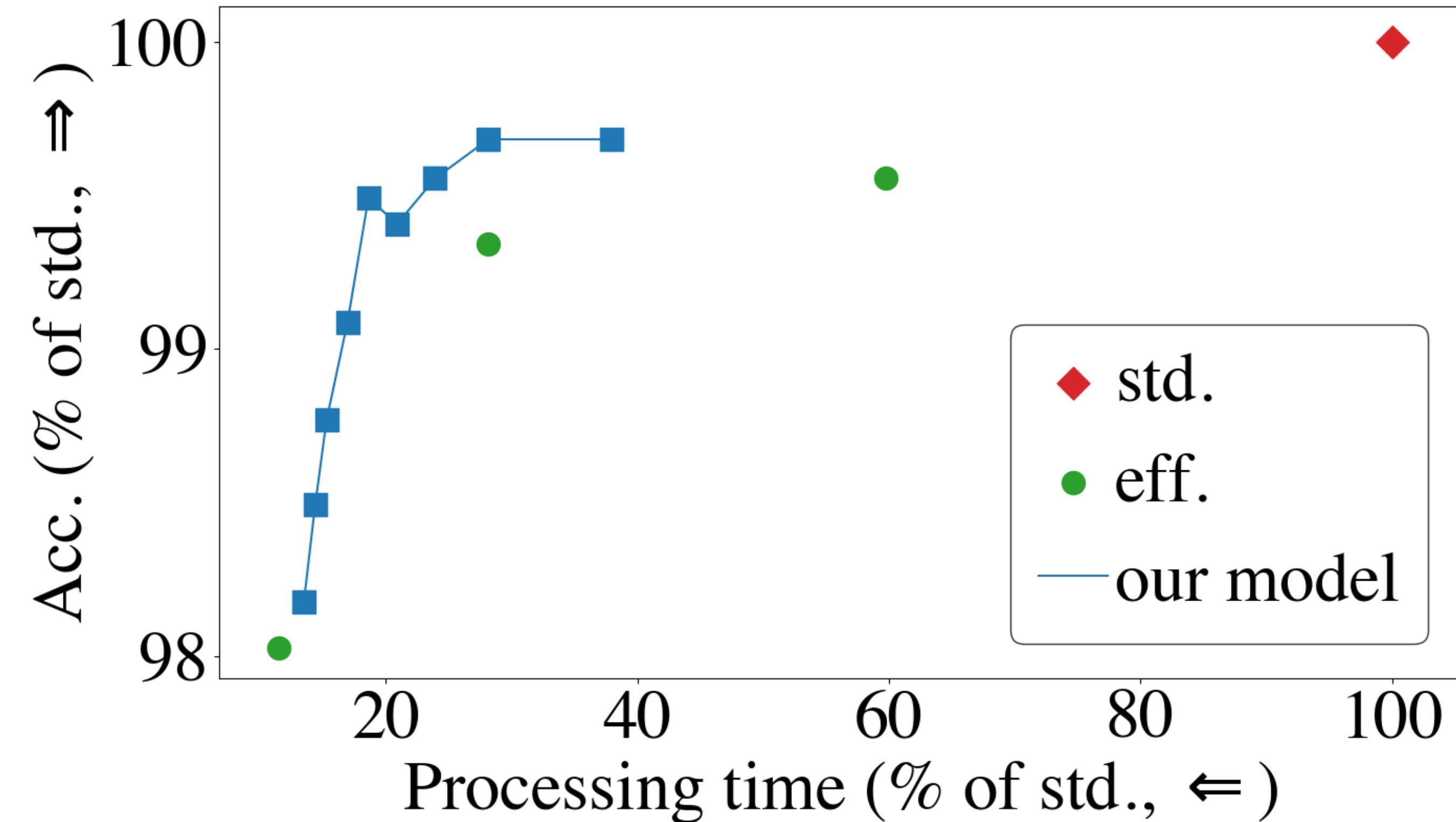


3 times faster, within 1% of full model



Better Speed/Accuracy Tradeoff

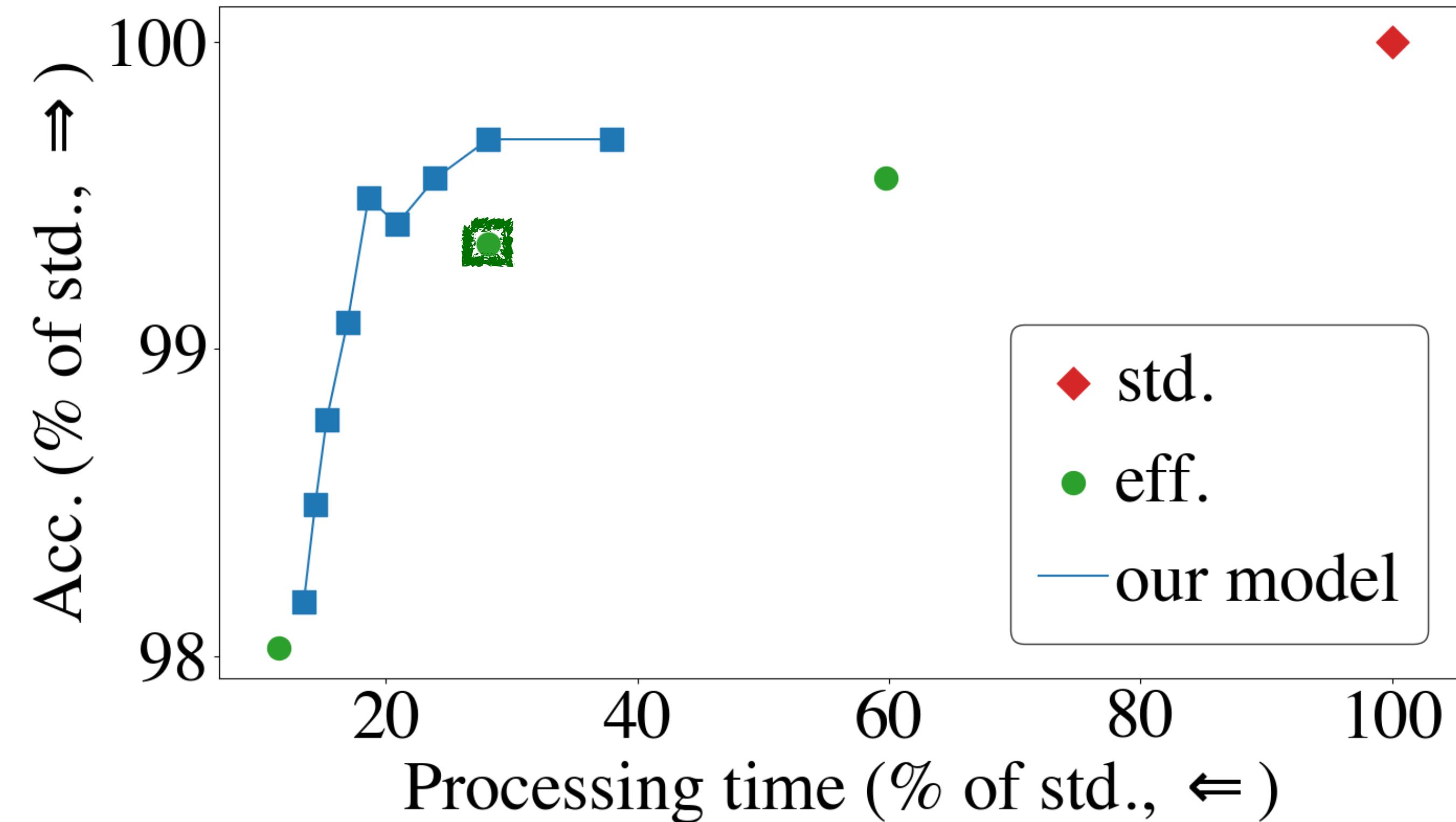
AG





Better Speed/Accuracy Tradeoff

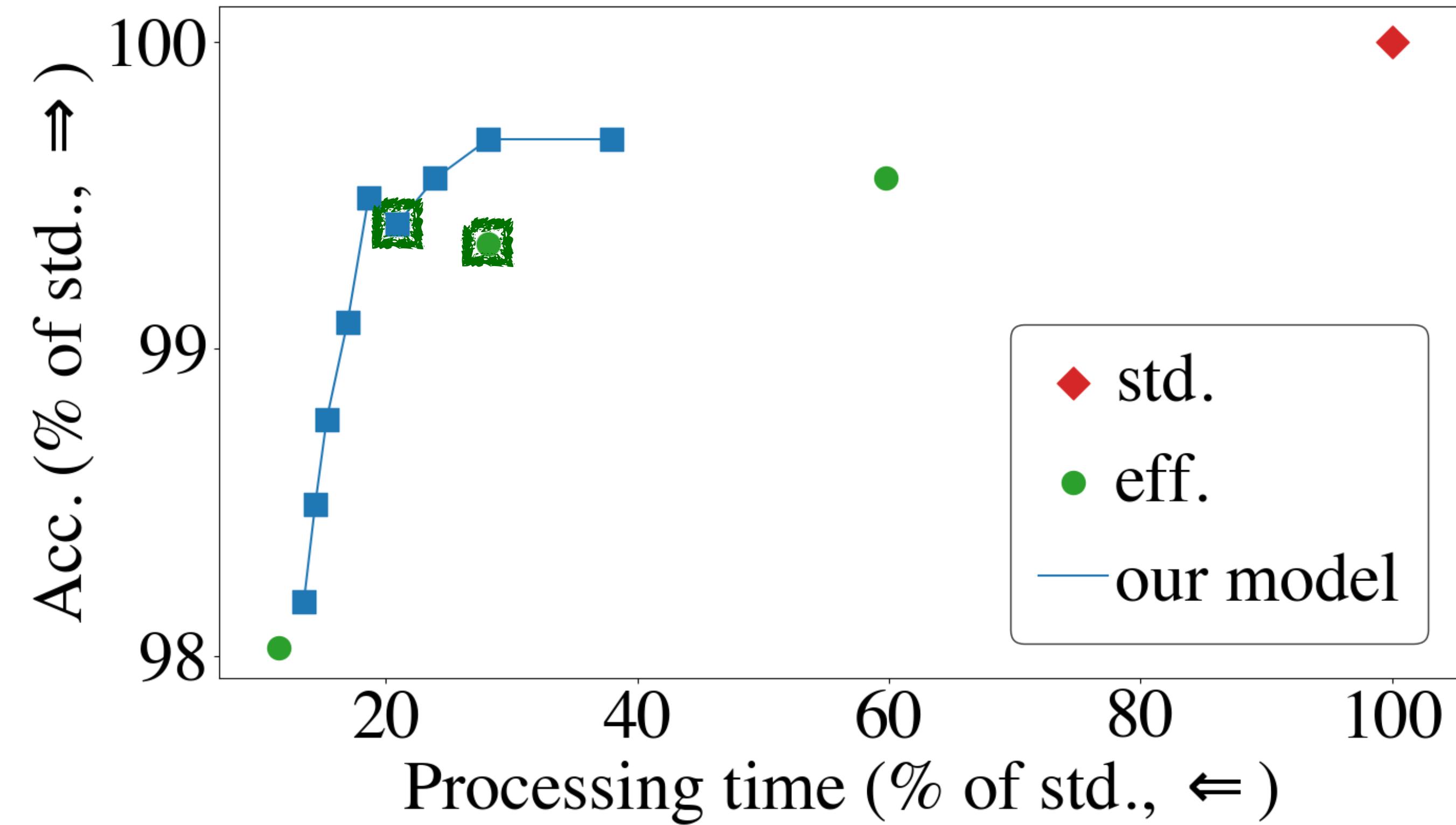
AG





Better Speed/Accuracy Tradeoff

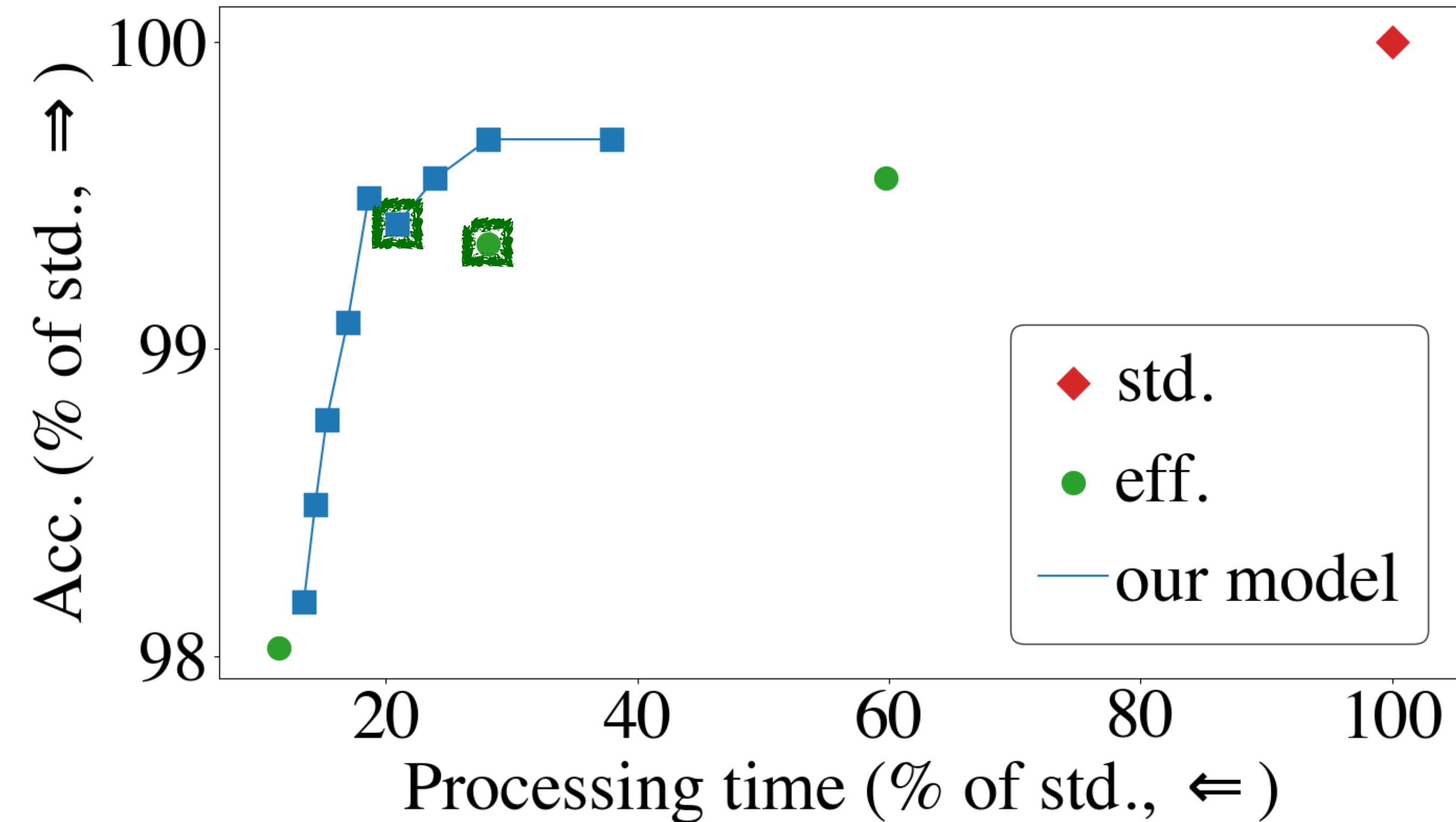
AG





Better Speed/Accuracy Tradeoff

AG



5 times faster, within 1% of full model



Efficiency

Open Questions

- What makes a good sparse structure?

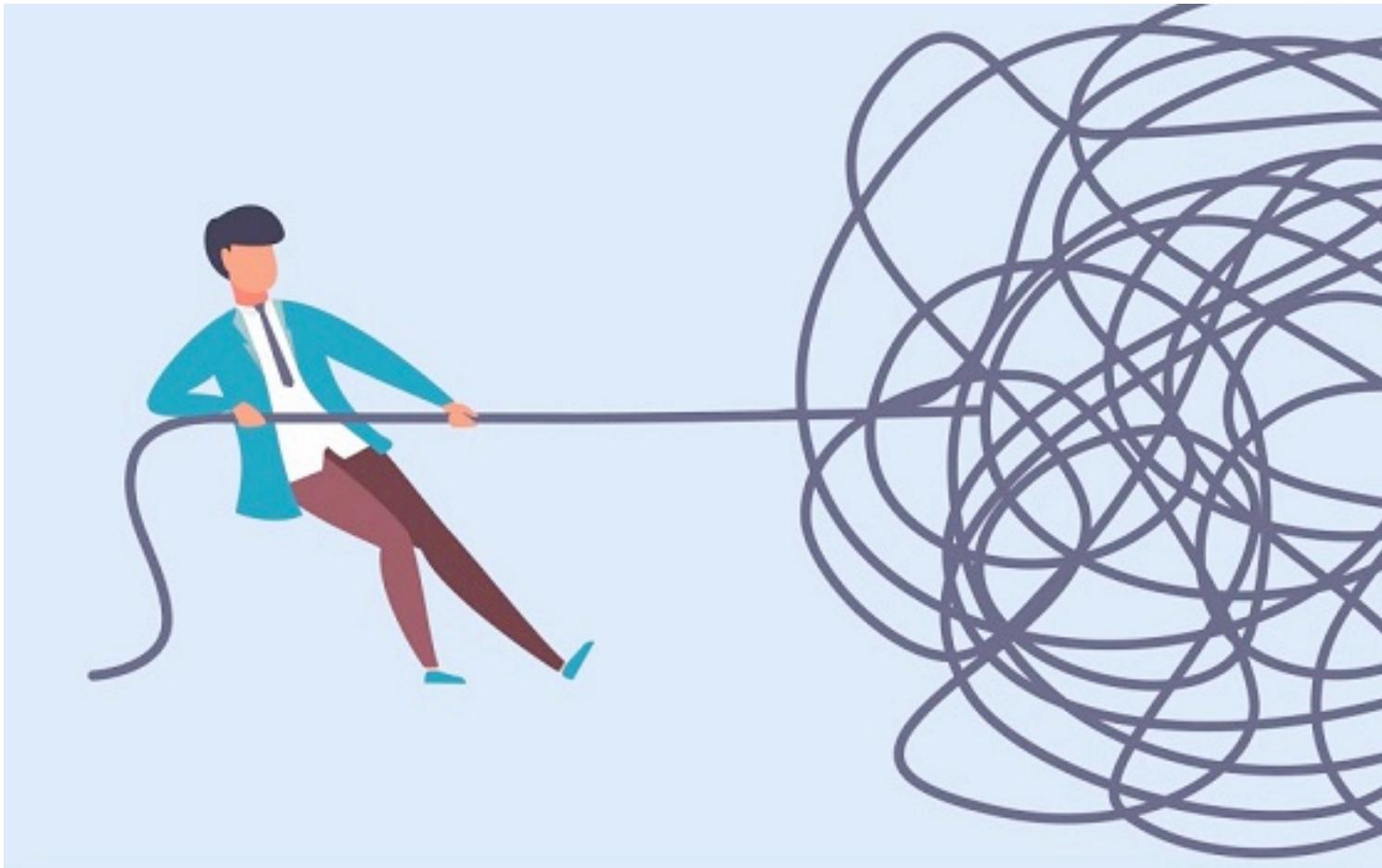


Efficiency

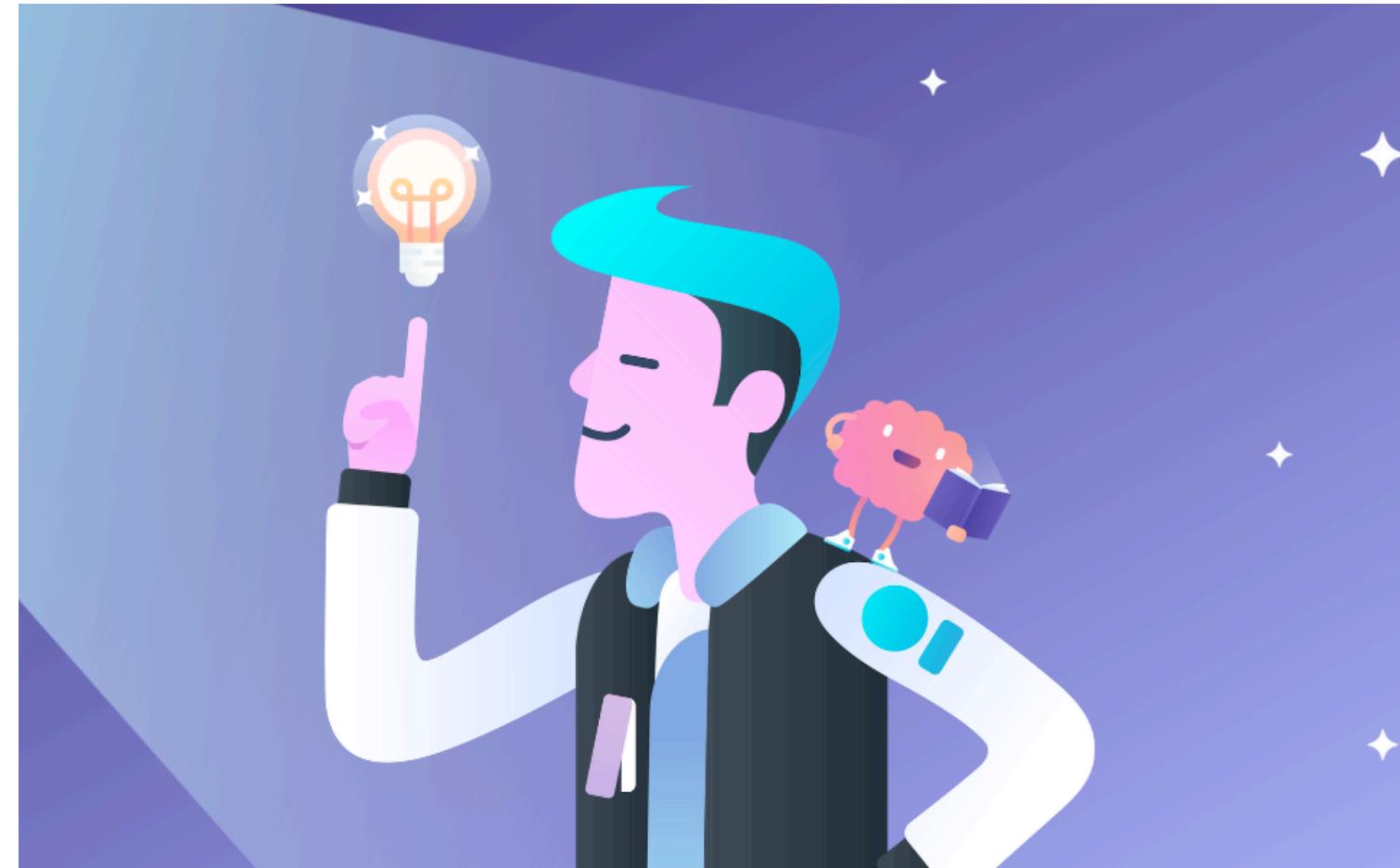
Open Questions

- What makes a good sparse structure?
- Combining different methods

Outline



Challenges

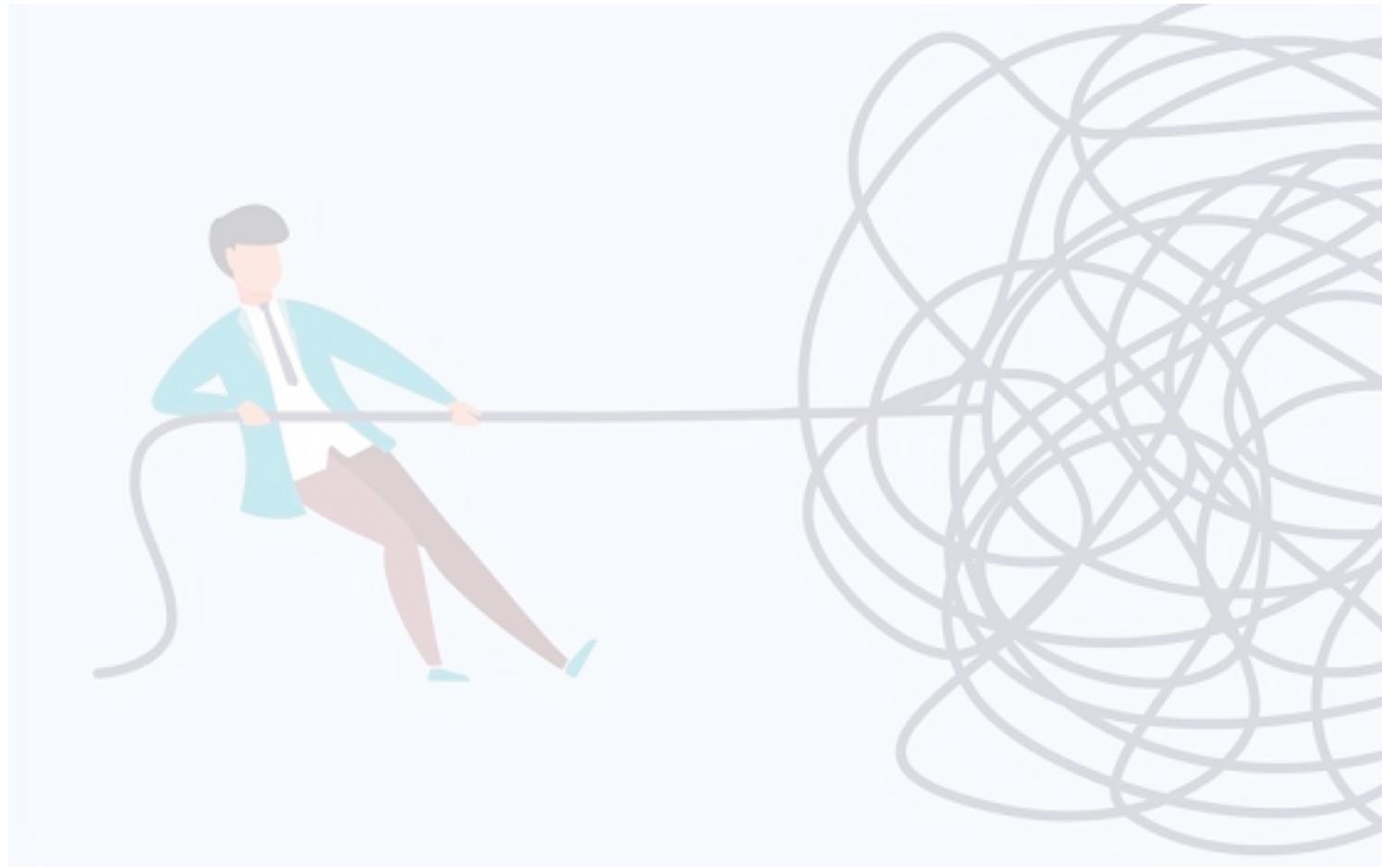


Opportunities



Mitigation

Outline



Challenges



Opportunities



Mitigation



Tackling Climate Change with Machine Learning

David Rolnick^{1*}, Priya L. Donti², Lynn H. Kaack³, Kelly Kochanski⁴, Alexandre Lacoste⁵, Kris Sankaran^{6,7}, Andrew Slavin Ross⁹, Nikola Milojevic-Dupont^{10,11}, Natasha Jaques¹², Anna Waldman-Brown¹², Alexandra Luccioni^{6,7}, Tegan Maharaj^{6,8}, Evan D. Sherwin², S. Karthik Mukkavilli^{6,7}, Konrad P. Körding¹, Carla Gomes¹³, Andrew Y. Ng¹⁴, Demis Hassabis¹⁵, John C. Platt¹⁶, Felix Creutzig^{10,11}, Jennifer Chayes¹⁷, Yoshua Bengio^{6,7}

¹University of Pennsylvania, ²Carnegie Mellon University, ³ETH Zürich, ⁴University of Colorado Boulder,

⁵Element AI, ⁶Mila, ⁷Université de Montréal, ⁸École Polytechnique de Montréal, ⁹Harvard University,

¹⁰Mercator Research Institute on Global Commons and Climate Change, ¹¹Technische Universität Berlin,

¹²Massachusetts Institute of Technology, ¹³Cornell University, ¹⁴Stanford University,

¹⁵DeepMind, ¹⁶Google AI, ¹⁷Microsoft Research

	Computer vision	NLP	Time-series analysis	Unsupervised learning	RL & Control	Causal inference	Uncertainty quantification	Transfer learning	Interpretable ML	Other
Electricity Systems	1	1.1	1.1 1.2	1	1.1		1.1 1.2	1.3	1.1	1.1
Transportation	2.1 2.2 2.4		2	2.1 2.4	2	2.1 2.4	2	2.1 2.4	2	
Buildings & Cities	3.2	3.3	3	3	3.1	3.1	3.3	3		
Industry	4.1 4.3		4.3	4.3	4	4.2 4.3		4.2 4.3	4.3	
Farms & Forests	5.1 5.3 5.4				5.2			5.4		
CO ₂ Removal			6.3				6.3	6.3		6.2
Climate Prediction	7.1		7				7.3		7	
Societal Impacts	8.1 8.4	8.4	8.2 8.3		8.2	8.3	8.2	8.1	8.3	
Solar Geoengineering			9.3		9.4		9.3 9.4			9.2
Tools for Individuals	10.1	10.1	10.2	10.3	10.2	10.1		10.2	10.2	
Tools for Society		11.1	11.2 11.1	11.3	11.2 11.1	11.1 11.3	11.1	11	11.1	11.1
Education		12.2			12.1					
Finance		13.2	13				13.2			



Claim Analysis

- Claim verification datasets
 - A claim is potentially verifiable if it is a) **well-formed** and b) **subjectively investigable**
 - CLIMATE-FEVER (Diggelmann et al., 2020)
- Claim verification tools
 - Wang et al. (2021)
 - Bhatia et al. (2021)
- Claim detection
 - Stammbach et al. (2022)



Document Analysis

- Topic detection
 - ClimaText (Varini et al., 2020)
- Climate report analysis
 - Bingler et al. (2021)
 - Friederich et al. (2021)
- Generation of evidence maps
 - Callaghan et al. (2021)



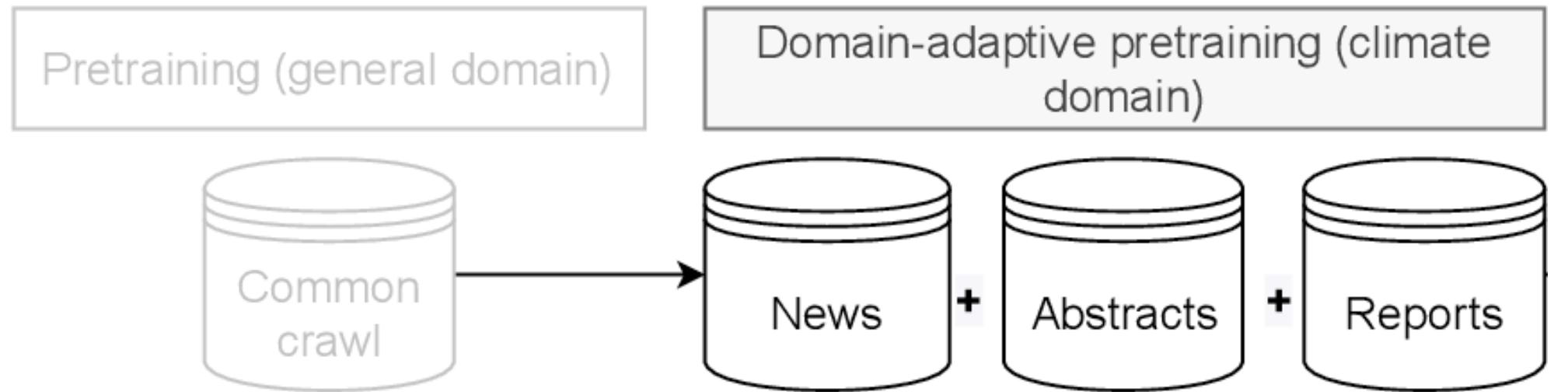
Social Media Analysis

- Sentiment, emotion and opinion analysis
 - Jiang et al. (2017)
 - Loureiro & Alló (2020)
 - El Barachi et al. (2021)
- Identification of neutralization and denial
 - Chen et al. (2019)
 - Bhatia et al. (2021)



General Purpose Tools

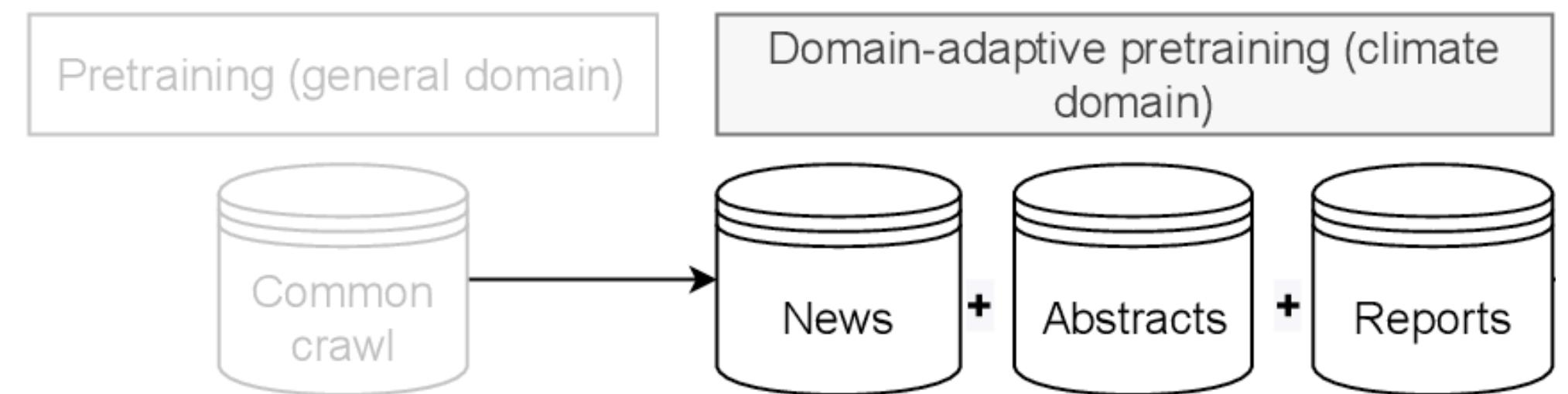
- ClimateBERT (Webersinke et al., 2022)
 - A pre-trained models trained on climate related texts





General Purpose Tools

- ClimateBERT (Webersinke et al., 2022)
 - A pre-trained models trained on climate related texts
- ClimateQA (Luccioni et al., 2020)
 - A question-answering dataset focusing on climate topics



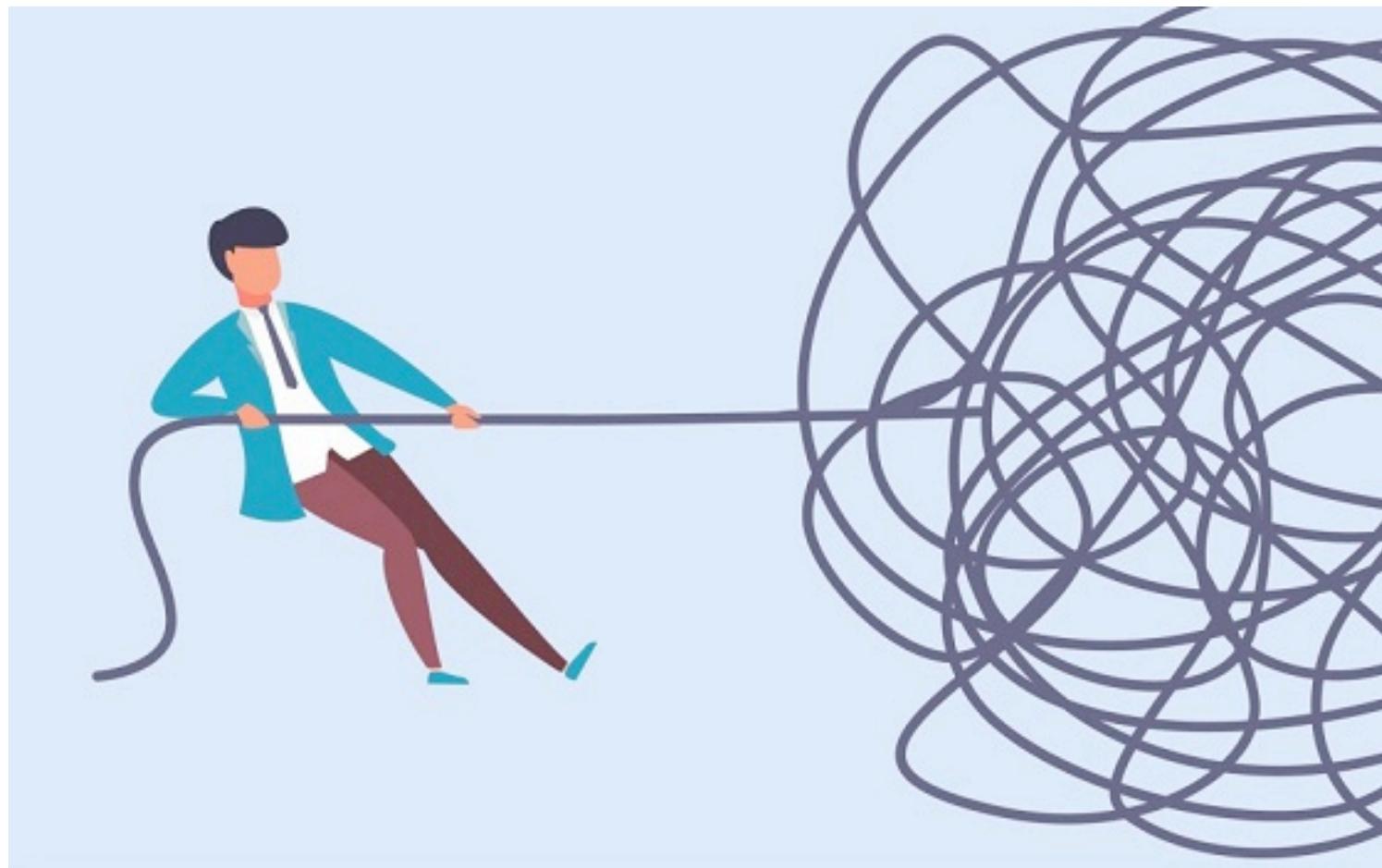
TCFD Question	Answer Passage
Does the organization describe the board's (or board committee's) oversight of climate-related risks and/or opportunities?	<i>The Company's Audit Committee has the delegated risk management oversight responsibility and receives updates on the risk management processes and key risk factors on a quarterly basis.</i>
Does the organization describe the climate-related risks or opportunities the organization has identified?	<i>The availability and price of these commodities are subject to factors such as changes in weather conditions, plantings, and government policies</i>

How can AI Inform Policy Makers?

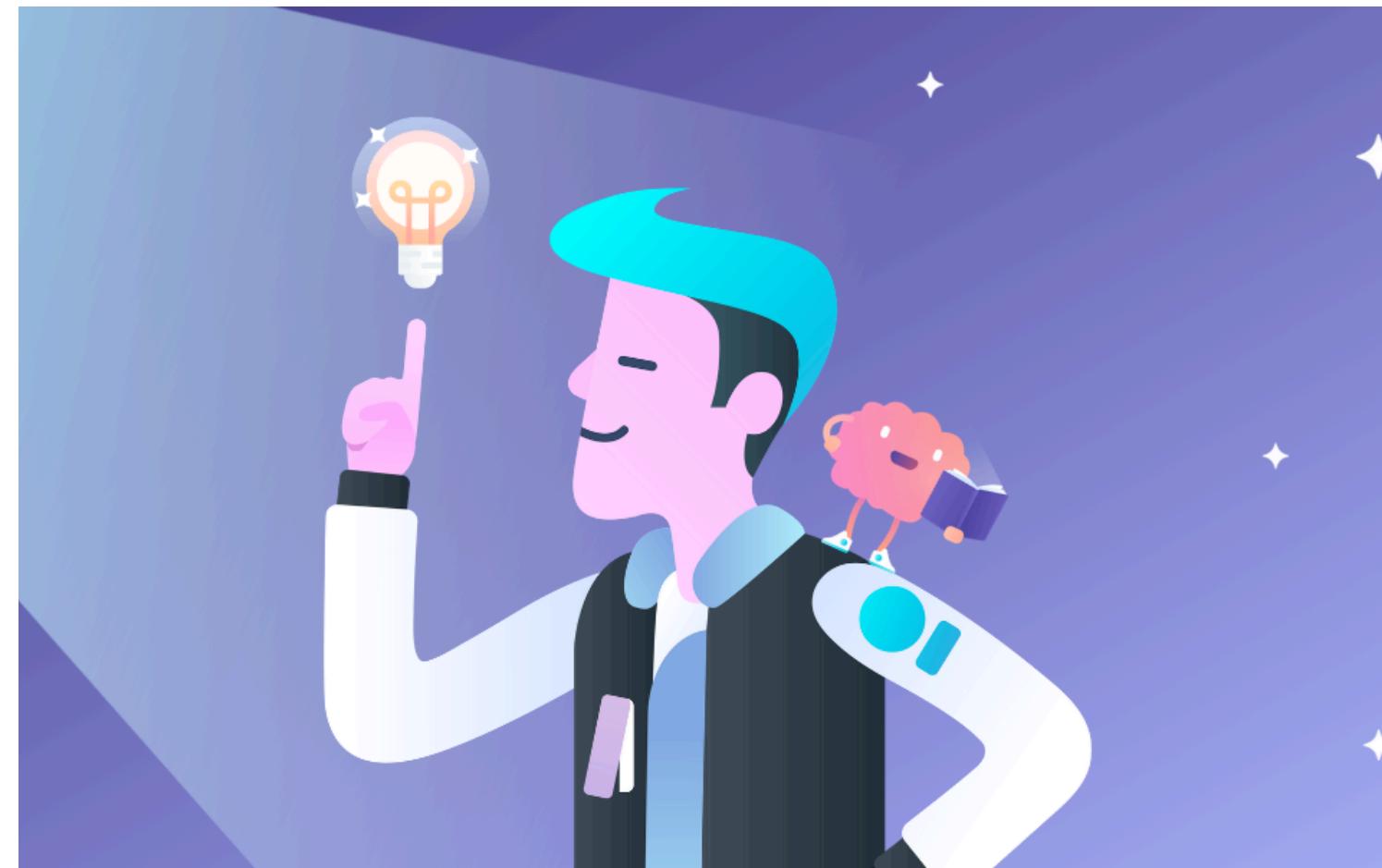


- Supporting policy-makers and activists
 - Stede & Patz (2021)
- A knowledge platform for holistic and effective climate action
 - Swarnakar & Modi (2021)

Summary



Challenges



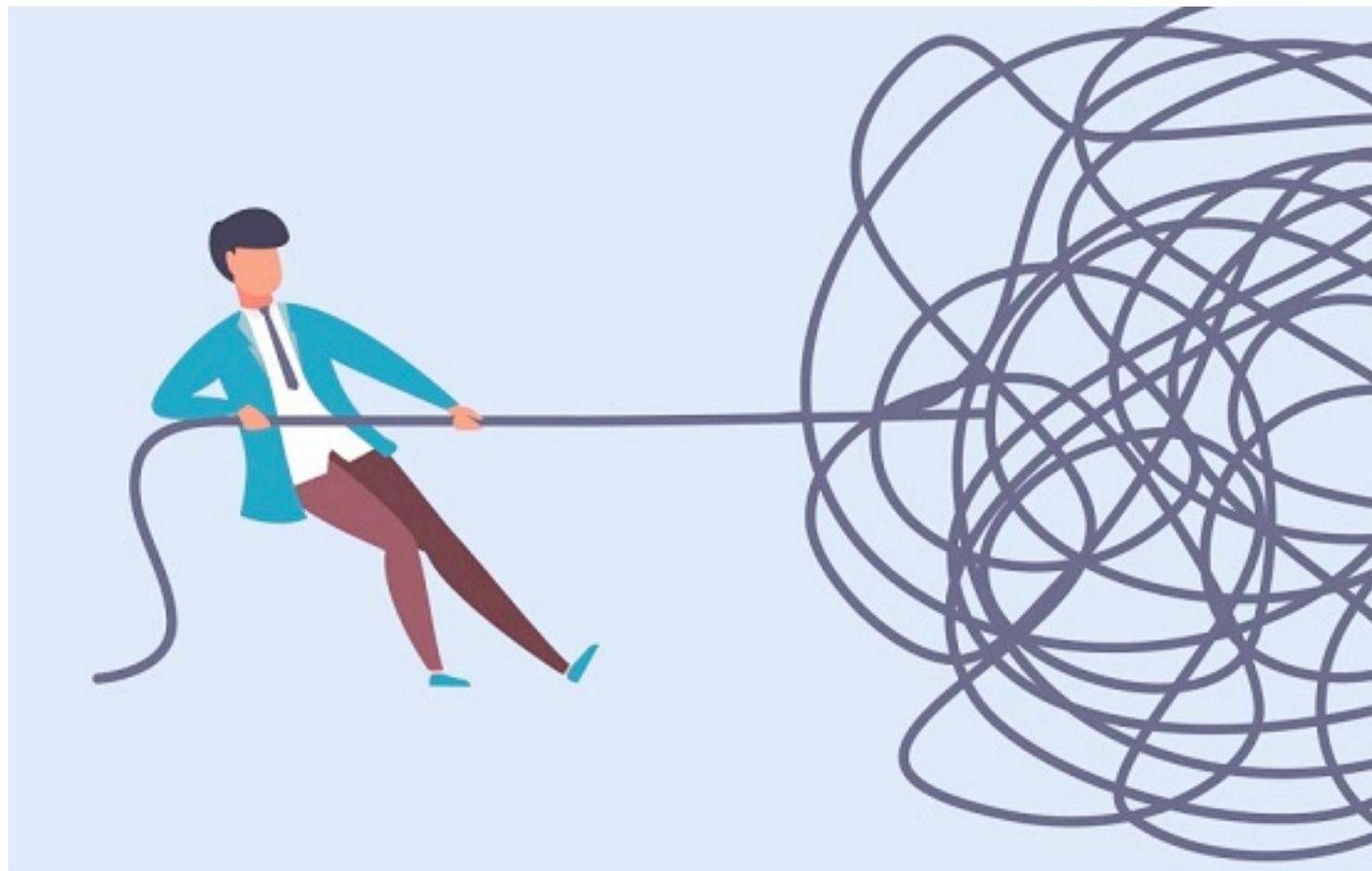
Opportunities



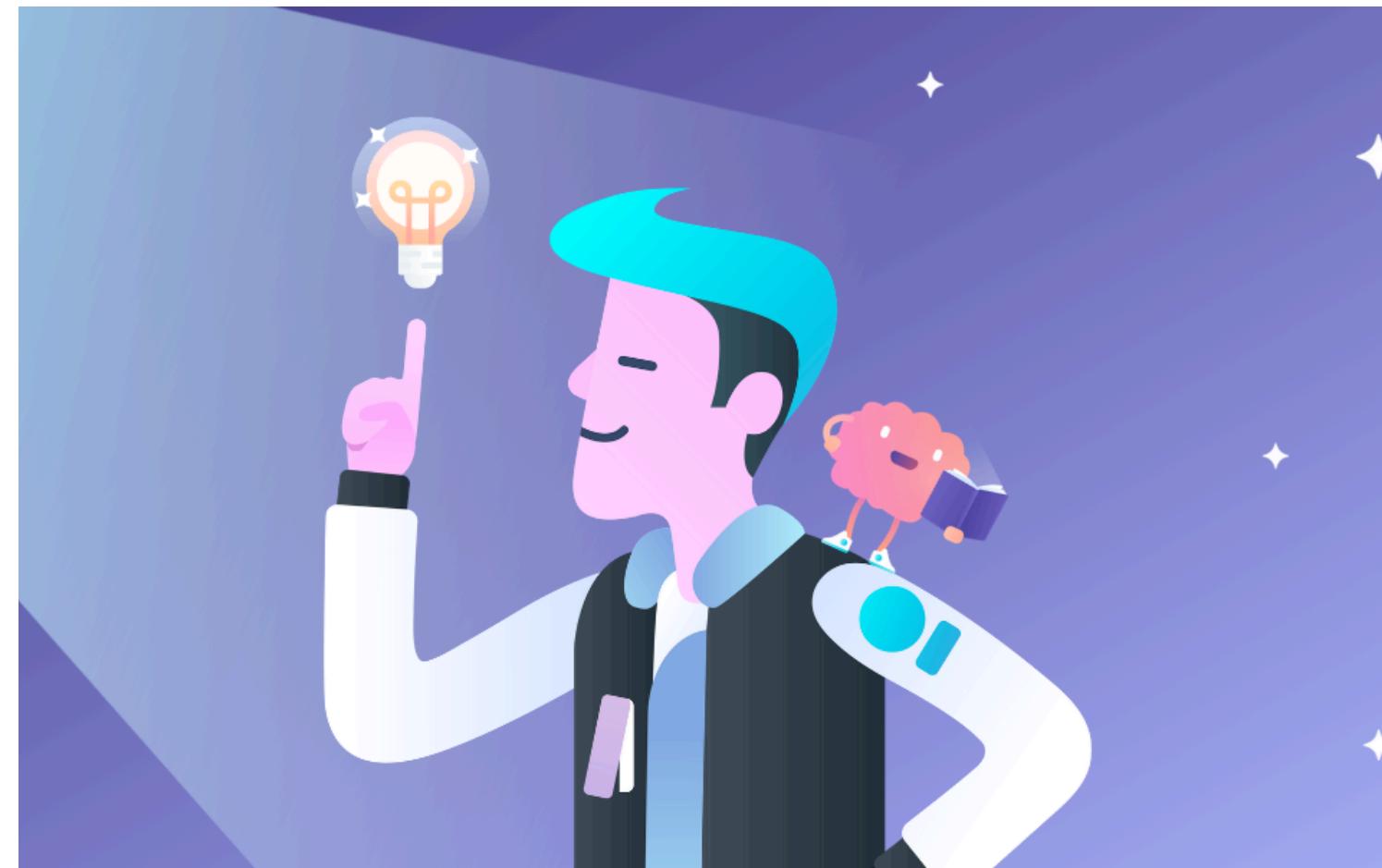
Mitigation

Thank you

Summary



Challenges



Opportunities



Mitigation